



Stefan, Thorsten (2016) Mathematical modelling of the selective forces maintaining diversity at the Major Histocompatibility Complex. PhD thesis, University of Glasgow.

<http://theses.gla.ac.uk/7029/>

Copyright and moral rights for this thesis are retained by the author

A copy can be downloaded for personal non-commercial research or study, without prior permission or charge

This thesis cannot be reproduced or quoted extensively from without first obtaining permission in writing from the Author

The content must not be changed in any way or sold commercially in any format or medium without the formal permission of the Author

When referring to this work, full bibliographic details including the author, title, awarding institution and date of the thesis must be given

MATHEMATICAL MODELLING OF THE
SELECTIVE FORCES MAINTAINING
DIVERSITY AT THE MAJOR
HISTOCOMPATIBILITY COMPLEX

THORSTEN STEFAN, M. RES.

SUBMITTED IN FULFILMENT OF THE REQUIREMENTS FOR THE DEGREE OF
Doctor of Philosophy

INSTITUTE OF BIODIVERSITY, ANIMAL HEALTH AND
COMPARATIVE MEDICINE

COLLEGE OF MEDICAL, VETERINARY AND LIFE SCIENCES
UNIVERSITY OF GLASGOW

JANUARY 2016

© THORSTEN STEFAN, M. RES.

Abstract

This work explores the long-standing question of which forces drive the maintenance of MHC diversity in vertebrates. More precisely, it investigates whether a special form of heterozygote advantage can explain the characteristic features of MHC genes, as the large number of alleles, the characteristic distribution of allele frequencies, and the trans-species polymorphism. This special overdominance variant is based on the divergent allele advantage hypothesis (Wakeland et al., 1990), and therefore the corresponding model is called the divergent allele advantage model. The novel tool of diversity profiles (${}^qD^Z$) will be used to characterise and compare allelic diversity of different overdominance models, and to contrast simulation results with observed data.

Chapters 2 and 3 perform a stability analysis of a n -allele system at equilibrium. Different overdominance models are compared in a common framework, and a novel model, based on the divergent allele advantage hypothesis, is introduced while its properties are analysed. Chapter 4 explores a more realistic scenario that includes mutation and drift and compares a modified divergent allele advantage model to traditional models of heterozygote advantage. Implications of the results of the previous chapters are discussed in chapter 5. Finally, chapter 6 extends the analysis of the divergent allele advantage model to a structured population, thus additionally including migration, and explores whether the findings of the previous chapters also hold in a more realistic scenario of a metapopulation of interacting subpopulations.

The results of this work demonstrate that the novel model based on the divergent allele advantage hypothesis can explain the main features of the diversity of the Major Histocompatibility Complex genetic region, but it predicts the existence of many genotypes that increase susceptibility to disease.

Table of Contents

Acknowledgements	20
Author's Declaration	20
1 General Introduction	21
1.1 Why are MHC genes so diverse?	21
1.2 Modelling the selective forces that maintain MHC polymorphism	26
1.3 Structure and function of the MHC	29
1.4 Measuring diversity	31
2 The Maintenance of MHC Alleles at Equilibrium – A General Framework	34
2.1 Introduction	34
2.2 Materials and Methods	35
2.2.1 General description of a single locus model	35
2.2.2 The genotype fitness matrix	37
2.3 Applications to selected overdominance models	38
2.3.1 Symmetric overdominance	38
2.3.2 Simple asymmetric overdominance	38
2.3.3 Reinforcing asymmetric overdominance	39
2.3.4 Divergent allele advantage	40
2.4 Discussion and Conclusions	42
3 MHC Diversity at Equilibrium under Different Overdominance Models	44
3.1 Introduction	44
3.2 Materials and Methods	45
3.2.1 Description of the allele frequency update algorithm	45

3.2.2	Allele configurations explored	45
3.2.3	Overdominance models and diversity measures	46
3.3	Results	47
3.3.1	Simple measures of persisting allelic diversity in the different overdominance models	47
3.3.2	Diversity profiles of selected overdominance models	52
3.3.3	Diversity in the DAA model	54
3.3.3.1	Number of persisting alleles and susceptible recognition sites of heterozygotes	54
3.3.3.2	Number of persisting alleles and population fitness	57
3.3.3.3	Diversity profiles for different sets of alleles in the DAA model	60
3.4	Discussion	63
3.5	Conclusions	66
3.5.1	Model limitations	67
3.5.2	Perspectives	67
4	The Maintenance of MHC Diversity in a Finite Population Subject to Evolutionary Forces	69
4.1	Preface	69
4.2	Introduction	69
4.3	Materials and Methods	71
4.3.1	Model overview	71
4.3.2	Changes in allele frequencies per generation	72
4.3.3	Selection under different overdominance models	73
4.3.4	Mutations	74
4.3.5	Parameter space explored by evolutionary simulations	75
4.3.5.1	Heterozygote advantage	78
4.3.6	Stochastic simulations	79
4.3.7	Comparison with MHC diversity in natural populations	79
4.4	Results	82

4.4.1	Model verification	82
4.4.2	Comparison of Overdominance Models	85
4.4.2.1	Number of alleles	85
4.4.2.2	Allele frequency distributions	87
4.4.2.3	Persisting variation in allele intrinsic merit	91
4.4.2.4	Diversity expressed as heterozygosity	95
4.4.2.5	Diversity expressed as gene pool share of the most frequent alleles	97
4.4.2.6	Diversity expressed as diversity profiles	97
4.4.2.7	Trans-species polymorphism	103
4.4.3	Sensitivity towards mutation rate and number of amino acids per allele	106
4.5	Discussion	112
4.6	Conclusion	113
5	MHC Diversity and Disease Susceptibility	114
5.1	Introduction	114
5.2	Materials and Methods	115
5.3	Results	115
5.3.1	Classifying alleles with low intrinsic merit	115
5.3.1.1	Allele intrinsic merit thresholds	115
5.3.1.2	Persistence times of alleles with low intrinsic merit . . .	117
5.3.1.3	Intermediate and transient alleles	117
5.3.1.4	Range of allele intrinsic merits	127
5.3.1.5	Weighted standard deviation in allele intrinsic merit . . .	127
5.3.1.6	Number of stable and intermediate alleles	130
5.3.2	Population fitness and alleles with low intrinsic merit	130
5.3.3	Population Fitness during and after a severe population bottleneck .	132
5.4	Discussion and Conclusions	134

6	MHC Diversity in Well-mixed and Structured Populations	137
6.1	Introduction	137
6.2	Materials and Methods	138
6.2.1	Split of the well-mixed population into subpopulations	138
6.2.2	Migration	138
6.2.3	Diversity profiles of the metapopulations	139
6.2.4	Overview of the simulations performed	140
6.3	Results	142
6.3.1	Number of alleles	142
6.3.2	Distribution of intrinsic merit of alleles	145
6.3.2.1	Range of allele intrinsic merits	147
6.3.2.2	Weighted standard deviation in allele intrinsic merit	147
6.3.3	Diversity profiles	148
6.3.3.1	Naive diversity	148
6.3.3.2	Similarity-sensitive diversity	152
6.3.3.3	Comparison between model results and observed allele frequencies	152
6.3.4	Trans-species polymorphism	159
6.4	Discussion and Conclusions	159
7	General Conclusions	161
A	Additional figures for chapter 3	163
B	Additional figures for chapter 4	166
C	Additional figures for chapter 5	233
D	Additional figures for chapter 6	252
	Bibliography	290

List of Tables

3.1	Initial sets of alleles in two different scenarios	46
3.2	Properties of of persisting alleles in different overdominance models	49
4.1	Parameter ranges and default values used for the simulations	76
4.2	Model parameters controlling selection for the 6 settings examined	77
4.3	Bovidae populations used for calculating measures of allelic diversity . . .	81
4.4	Comparison of simulated heterozygosities to theoretical values	82
5.1	Population bottlenecks of four different severity levels for the DAA and SAsOD models	132
6.1	Schedule of simulations performed for the sensitivity analysis of migration rate and interval	141
6.2	Schedule of simulations performed for the sensitivity analysis of the number of subpopulations	142

List of Figures

2.1	Genotype fitness in the DAA model based on epitopes recognised	41
3.1	Distribution of allele frequencies in the SOD model	47
3.2	Distribution of allele frequencies in the SAsOD model	48
3.3	Distribution of allele frequencies in the RAsOD model	48
3.4	Illustrative distribution of allele frequencies in the DAA model	49
3.5	Number of alleles maintained in the DAA model	50
3.6	Range of intrinsic merits of alleles maintained in the DAA model	51
3.7	Combined diversity measure of alleles at equilibrium in the DAA model . .	51
3.8	Diversity profiles for alleles maintained in the DAA model	53
3.9	Correlation between number of alleles and allele intrinsic merit range . . .	54
3.10	Susceptible recognition sites vs. number of alleles maintained (taking all alleles into account)	56
3.11	Susceptible recognition sites vs. number of alleles maintained (taking only the most frequent alleles into account)	57
3.12	Susceptible recognition sites vs. number of alleles maintained (taking only a variable number of the most frequent alleles into account)	58
3.13	Population fitness vs. number of alleles maintained	59
3.14	Population fitness vs. susceptible recognition sites	59
3.15	Diversity profiles for naive diversity	61
3.16	Diversity profiles for similarity-sensitive diversity based on intrinsic merit differences	62
3.17	Diversity profiles for similarity-sensitive diversity based on recognition site differences	62
3.18	Susceptible recognition sites vs. recognition site differences	64

4.1	Positions of the 6 settings examined	77
4.2	Heterozygote advantage for different overdominance models	79
4.3	Average sequence difference, heterozygote advantage and population fitness for different overdominance models	80
4.4	Comparison of simulated to theoretical frequency spectra for neutral alleles	83
4.5	Comparison of simulated to theoretical frequency spectra for the FSOD model	84
4.6	Number of alleles over time in the DAA model	85
4.7	Number of alleles vs. population size in different overdominance models . .	86
4.8	Population size and number of distinct alleles in a population (DAA model)	87
4.9	Number of distinct alleles in a population (DAA model) when varying two main parameters	88
4.10	Weighted frequency spectra for different models	89
4.11	Variation in allele intrinsic merit in the DAA and SAsOD models	92
4.12	Standard deviation in allele intrinsic merit in the DAA and SAsOD models .	92
4.13	Share of alleles with low intrinsic merit for the DAA and SAsOD models .	93
4.14	Mean age of alleles with low intrinsic merit for the DAA and SAsOD models	94
4.15	Comparison of predictions from different overdominance models to expected heterozygosities calculated from published allele frequencies	96
4.16	Comparison of predictions from different overdominance models to the share of the 5 most frequent alleles calculated from published allele frequencies .	98
4.17	Diversity profiles (naive diversity) for different overdominance models . . .	99
4.18	Diversity profiles (naive diversity) for the DAA and the SAsOD models . .	100
4.19	Diversity profiles based on allele intrinsic merit differences for different overdominance models	101
4.20	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models	102
4.21	Diversity profiles based on differences of encoded amino acid sequences for different overdominance models	103
4.22	Diversity profiles based on differences of encoded amino acid sequences for the DAA and the SAsOD models	104
4.23	Diversity profiles comparing observed data to simulation results from the DAA and SAsOD models	105

4.24	Distribution of emergence times of alleles for different overdominance models	107
4.25	Sensitivity of the number of alleles when varying mutation rate and number of variable sites	110
4.26	Sensitivity of the distribution of allele intrinsic merits when varying mutation rate and number of variable sites	110
4.27	Sensitivity of the distribution of emergence times of alleles when varying mutation rate and number of variable sites	111
4.28	Sensitivity of the expected heterozygosity when varying mutation rate and number of variable sites	111
4.29	Sensitivity of diversity profiles when varying mutation rate and number of variable sites	112
5.1	Allele intrinsic merit range when taking all alleles into account	118
5.2	Correlation between allele intrinsic merit, allele frequency and time of allele emergence in the DAA and SAsOD models	119
5.3	Number of intermediate and transient alleles in the DAA and SAsOD models using the less restrictive intrinsic merit threshold	120
5.4	Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models using the less restrictive intrinsic merit threshold	121
5.5	Frequency of intermediate alleles in the DAA and SAsOD models (stricter allele intrinsic merit threshold)	123
5.6	Frequency of intermediate alleles in the DAA and SAsOD models (less restrictive allele intrinsic merit threshold)	124
5.7	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model	126
5.8	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 1	128
5.9	Allele intrinsic merit range after transient alleles were removed	129
5.10	Standard deviation in allele intrinsic merit after transient alleles were removed	129
5.11	Number of alleles after transient alleles were removed	130
5.12	Population fitness before and after removal of alleles with low intrinsic merit for the DAA and SAsOD models	131
5.13	Population fitness before and after a severe population bottleneck for the DAA and SAsOD models	133

5.14	Population fitness vs. time during a severe population bottleneck for the DAA and SAsOD models	134
6.1	Number of alleles over time for well-mixed and structured populations . . .	143
6.2	Final number of alleles for subpopulations with different migration rates . .	144
6.3	Final number of alleles for different numbers of subpopulations in a star arrangement (1)	145
6.4	Distribution of intrinsic merit of alleles for the DAA and asymmetric overdominance models when varying the migration rate	146
6.5	Diversity profiles (naive diversity) when varying the migration rate	149
6.6	Diversity profiles (naive diversity) when varying the number of subpopulations in a star arrangement and the migration interval	150
6.7	Diversity profiles (naive diversity) when varying the arrangement of subpopulations	151
6.8	Diversity profiles when accounting for allele intrinsic merit differences and varying the migration rate and the number of subpopulations	153
6.9	Diversity profiles when accounting for allele sequence divergence and varying migration rate and number of subpopulations	154
6.10	Expected heterozygosity for predictions from the DAA and asymmetric overdominance models and observed data	156
6.11	Share of the five most frequent alleles for predictions from the DAA and asymmetric overdominance models and observed data	157
6.12	Diversity profiles (naive diversity) for the DAA and asymmetric overdominance models	158
A.1	Proportion of susceptible recognition sites vs. allele intrinsic merit range . .	164
A.2	Population fitness vs. allele intrinsic merit range	165
B.1	Heterozygote advantage for different overdominance models and settings 3, 4, 5 and 6	167
B.2	Average sequence difference, heterozygote advantage and population fitness for the DAA model and a population size of 100000	168
B.3	Average sequence difference, heterozygote advantage and population fitness for the DAA model and a population size of 10 million	169
B.4	Number of alleles over time in the DAA model for settings 3, 4, 5 and 6 . .	170

B.5	Number of alleles vs. population size for different overdominance models and settings 3, 4, 5 and 6	171
B.6	Weighted frequency spectra for the SOD and SAsOD models and settings 2 and 5	172
B.7	Weighted frequency spectra for the DAA and SAsOD models and settings 3, 4, 5 and 6	173
B.8	Weighted frequency spectra for the DAA and SAsOD models, setting 1 and different population sizes	174
B.9	Weighted frequency spectra for the DAA and SAsOD models, setting 3 and different population sizes	175
B.10	Weighted frequency spectra for the DAA and SAsOD models, setting 5 and different population sizes	176
B.11	Variation in allele intrinsic merit for the DAA and SAsOD models and settings 3, 4, 5 and 6	177
B.12	Variation in allele intrinsic merit for the DAA and RAsOD models	178
B.13	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 1 and different population sizes	179
B.14	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 2 and different population sizes	180
B.15	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 3 and different population sizes	181
B.16	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 4 and different population sizes	182
B.17	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 5 and different population sizes	183
B.18	Variation in allele intrinsic merit for the DAA and SAsOD models, setting 6 and different population sizes	184
B.19	Standard deviation in allele intrinsic merit for the DAA and SAsOD models and settings 3, 4, 5 and 6	185
B.20	Share of alleles with low intrinsic merit for the DAA and SAsOD models and settings 3 and 4	186
B.21	Share of alleles with low intrinsic merit for the DAA and SAsOD models and settings 5 and 6	187

B.22	Share of alleles with low intrinsic merit for the SAsOD model and settings 1 and 2	188
B.23	Share of alleles with low intrinsic merit for the SAsOD and RAsOD models and settings 2 and 3	188
B.24	Share of alleles with low intrinsic merit for the SAsOD model and settings 3, 4, 5 and 6	189
B.25	Mean age of alleles with low intrinsic merit for the DAA and SAsOD models and settings 3 and 4	190
B.26	Mean age of alleles with low intrinsic merit for the DAA and SAsOD models and settings 5 and 6	191
B.27	Mean age of alleles with low intrinsic merit for the SAsOD model and settings 1 and 2	192
B.28	Mean age of alleles with low intrinsic merit for the SAsOD and RAsOD models and settings 2 and 3	192
B.29	Mean age of alleles with low intrinsic merit for the SAsOD model and settings 3, 4, 5 and 6	193
B.30	Comparison of predictions from the DAA and SAsOD models (settings 3, 4, 5 and 6) to expected heterozygosities calculated from published allele frequencies	194
B.31	Comparison of predictions from the DAA and SAsOD models (settings 3, 4, 5 and 6) to the share of the 5 most frequent alleles calculated from published allele frequencies	195
B.32	Diversity profiles (naive diversity) for different overdominance models and settings 3, 4, 5 and 6	196
B.33	Diversity profiles (naive diversity) for the DAA and SAsOD models and setting 2	197
B.34	Diversity profiles (naive diversity) for the DAA and SAsOD models and setting 3	198
B.35	Diversity profiles (naive diversity) for the DAA and SAsOD models and setting 4	199
B.36	Diversity profiles (naive diversity) for the DAA and SAsOD models and setting 5	200
B.37	Diversity profiles (naive diversity) for the DAA and SAsOD models and setting 6	201

B.38	Diversity profiles based on allele intrinsic merit differences for different overdominance models and settings 3, 4, 5 and 6	202
B.39	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models and setting 2	203
B.40	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models and setting 3	204
B.41	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models and setting 4	205
B.42	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models and setting 5	206
B.43	Diversity profiles based on allele intrinsic merit differences for the DAA and the SAsOD models and setting 6	207
B.44	Diversity profiles based on differences between encoded amino acid sequences for different overdominance models and settings 3, 4, 5 and 6	208
B.45	Diversity profiles based on differences between encoded amino acid sequences for the DAA and the SAsOD models and setting 2	209
B.46	Diversity profiles based on differences between encoded amino acid sequences for the DAA and the SAsOD models and setting 3	210
B.47	Diversity profiles based on differences between encoded amino acid sequences for the DAA and the SAsOD models and setting 4	211
B.48	Diversity profiles based on differences between encoded amino acid sequences for the DAA and the SAsOD models and setting 5	212
B.49	Diversity profiles based on differences between encoded amino acid sequences for the DAA and the SAsOD models and setting 6	213
B.50	Diversity profiles comparing observed data to predictions from the DAA and SAsOD models and settings 3, 4, 5 and 6	214
B.51	Diversity profiles comparing observed data to predictions from the SOD model and setting 5	215
B.52	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 1	216
B.53	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 2	217
B.54	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 3	218

B.55	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 4	219
B.56	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 5	220
B.57	Distribution of emergence times of alleles for the DAA and SAsOD models and setting 6	221
B.58	Distribution of emergence times of alleles for the RAsOD model and setting 3	222
B.59	Distribution of emergence times of alleles for the SOD model and setting 5	222
B.60	Sensitivity of the number of alleles in setting 2 when varying mutation rate and number of variable sites	223
B.61	Sensitivity of the range of allele intrinsic merits in settings 1 and 2 when varying mutation rate and number of variable sites	224
B.62	Sensitivity of the weighted standard deviation in allele intrinsic merit in settings 1 and 2 when varying mutation rate and number of variable sites . . .	225
B.63	Sensitivity of the distribution of allele intrinsic merits for setting 1 and a population size of 10000 when varying mutation rate and number of variable sites	226
B.64	Sensitivity of the distribution of allele intrinsic merits for setting 1 and population sizes of 1 million and 10 million when varying mutation rate and number of variable sites	227
B.65	Sensitivity of the distribution of allele intrinsic merits for setting 2 and population sizes of 10000 and 100000 when varying mutation rate and number of variable sites	228
B.66	Sensitivity of the distribution of allele intrinsic merits for setting 2 and a population size of 1 million when varying mutation rate and number of variable sites	229
B.67	Sensitivity of the distribution of emergence times of alleles for setting 1 and a population size of 10000 when varying mutation rate and number of variable sites	229
B.68	Sensitivity of the distribution of emergence times of alleles for setting 1 and population sizes of 1 million and 10 million when varying mutation rate and number of variable sites	230
B.69	Sensitivity of the distribution of emergence times of alleles for setting 2 and population sizes of 10000 and 100000 when varying mutation rate and number of variable sites	231

B.70	Sensitivity of the distribution of emergence times of alleles for setting 2 and a population size of 1 million when varying mutation rate and number of variable sites	232
C.1	Allele intrinsic merit range when taking all alleles into account for settings 3, 4, 5 and 6	234
C.2	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 2	235
C.3	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 3	236
C.4	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 4	237
C.5	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 5	238
C.6	Allele intrinsic merit range for different age thresholds that removed transient alleles and setting 6	239
C.7	Standard deviation in allele intrinsic merit after transient alleles were removed for settings 3, 4, 5 and 6	240
C.8	Number of alleles after transient alleles were removed for settings 3, 4, 5 and 6	241
C.9	Share of the gene pool of intermediate and transient alleles, using the less restrictive intrinsic merit threshold for settings 3 and 4	242
C.10	Share of the gene pool of intermediate and transient alleles, using the less restrictive intrinsic merit threshold and settings 5 and 6	243
C.11	Share of the gene pool of intermediate and transient alleles, using the stricter intrinsic merit threshold and settings 1 and 2	244
C.12	Share of the gene pool of intermediate and transient alleles, using the stricter intrinsic merit threshold and settings 3 and 4	245
C.13	Share of the gene pool of intermediate and transient alleles, using the stricter intrinsic merit threshold and settings 5 and 6	246
C.14	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model, using the less restrictive intrinsic merit threshold and settings 3 and 4	247
C.15	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model, using the less restrictive intrinsic merit threshold and settings 5 and 6	248

C.16	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model, using the stricter intrinsic merit threshold and settings 1 and 2	249
C.17	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model, using the stricter intrinsic merit threshold and settings 3 and 4	250
C.18	Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model, using the stricter intrinsic merit threshold and settings 5 and 6	251
D.1	Final number of alleles for different numbers of subpopulations in a star arrangement (2)	252
D.2	Range of intrinsic merits of alleles for different population sizes and different migration rates	253
D.3	Range of intrinsic merits of alleles for different numbers of subpopulations and different migration rates	254
D.4	Standard deviation in allele intrinsic merit for different population sizes and different migration rates	255
D.5	Standard deviation in allele intrinsic merit for different numbers of subpopulations and different migration rates	256
D.6	Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every generation	257
D.7	Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every 100 generations . . .	258
D.8	Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every 100000 generations .	259
D.9	Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement and no migration between subpopulations	260
D.10	Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every generation	261
D.11	Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every 100 generations . . .	262
D.12	Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every 10000 generations . .	263

D.13	Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement and no migration between subpopulations	264
D.14	Diversity profiles (naive diversity) for metapopulations of different connectivities and different migration rates	265
D.15	Diversity profiles (naive diversity) for metapopulations of different connectivities and a migration rate of 2 individuals every 100 generations	266
D.16	Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations, the DAA and asymmetric overdominance models and a migration rate of 2 individuals every generation	267
D.17	Diversity profiles (naive diversity) for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 100 generations	268
D.18	Diversity profiles (naive diversity) for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 10000 generations	269
D.19	Diversity profiles (naive diversity) for the DAA and asymmetric overdominance models and no migration between subpopulations	270
D.20	Diversity profiles based on allele intrinsic merit differences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every generation	271
D.21	Diversity profiles based on allele intrinsic merit differences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 100 generations	272
D.22	Diversity profiles based on allele intrinsic merit differences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 10000 generations	273
D.23	Diversity profiles based on allele intrinsic merit differences for the DAA and asymmetric overdominance models and no migration between subpopulations	274
D.24	Diversity profiles based on differences between encoded amino acid sequences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every generation	275
D.25	Diversity profiles based on differences between encoded amino acid sequences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 100 generations	276
D.26	Diversity profiles based on differences between encoded amino acid sequences for the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 10000 generations	277

D.27 Diversity profiles based on differences between encoded amino acid sequences for the DAA and asymmetric overdominance models and no migration between subpopulations	278
D.28 Expected heterozygosity for metapopulations with different connectivities and observed data	279
D.29 Share of the five most frequent alleles for metapopulations with different connectivities and observed data	280
D.30 Diversity profiles (naive diversity) for metapopulations with different connectivities of subpopulations	281
D.31 Distribution of emergence times of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every generation	282
D.32 Distribution of emergence times of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every 100 generations	283
D.33 Distribution of emergence times of alleles for a metapopulation in a star arrangement and a migration rate of 2 individuals every 100000 generations	284
D.34 Distribution of emergence times of alleles for a metapopulation in a star arrangement and no migration between subpopulations	285
D.35 Distribution of emergence times of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every generation	286
D.36 Distribution of emergence times of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every 100 generations	287
D.37 Distribution of emergence times of alleles for a metapopulation in a linear arrangement and a migration rate of 2 individuals every 100000 generations	288
D.38 Distribution of emergence times of alleles for a metapopulation in a linear arrangement and no migration between subpopulations	289

Acknowledgements

This work was supported by the EU funded Marie Curie Initial Training Network (MC-ITN) program.

I would further like to thank my two supervisors, Dr Louise Matthews and Professor Michael Stear for their continuous and incredibly valuable support in every respect throughout my time in Glasgow and beyond.

Author's Declaration

I hereby declare that I am the sole author of this thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners.

Chapter 1

General Introduction

1.1 Why are MHC genes so diverse?

Major histocompatibility (MHC) molecules play a central role in the immune system of vertebrates by presenting peptides to T-cells (Hughes and Yeager, 1998), hence conferring resistance or causing susceptibility to pathogens. Due to its importance in human medicine, detailed information about the MHC and its function is available (Hedrick, 1994).

The genes encoding the MHC are known to be the most polymorphic loci in vertebrates (Hedrick, 1994). Furthermore, alleles at these loci often differ at multiple nucleotides (Takahata and Nei, 1990; Hughes and Yeager, 1998; Stear et al., 2005). The rate of nonsynonymous substitutions in the peptide binding region (PBR, see section 1.3) is significantly higher than that of synonymous substitutions, while in other regions the nonsynonymous substitution rate is significantly lower (Hughes and Nei, 1988; Klein et al., 1993), as substitutions that change the amino acid composition are in most cases deleterious. Finally, many MHC polymorphisms are quite ancient, predating speciation events (Figuerola et al., 1988), and are therefore termed “trans-species” polymorphisms. For example, certain MHC alleles from human and chimpanzee belong to allelic lineages that existed before divergence of these two species (Lawlor et al., 1988; Mayer et al., 1988; Gyllenstein and Erlich, 1989). However, an alternative explanation for similarities between MHC alleles of related species has recently been proposed by Grossen et al. (2014), who demonstrated that introgression from domestic goats into Alpine ibex could likewise account for these allelic similarities.

All these observations present interesting questions from an evolutionary perspective. Without selection acting on the genes of the MHC, genetic diversity could be expected to be low and mainly controlled by mutation and genetic drift (Apanius et al., 1997). Such a neutral model would result in a low number of MHC alleles and strongly uneven allele frequencies at equilibrium (Apanius et al., 1997; Spurgin and Richardson, 2010). However, the observed

allele frequencies are much more even, as MHC polymorphisms are characterized by a large number of alleles occurring at intermediate frequencies (Meyer and Thomson, 2001) – a pattern of polymorphism inconsistent with selective neutrality but rather suggesting some form of balancing selection acting on MHC genes (Hedrick and Thomson, 1983).

The exact nature of these biological forces, capable of driving the extraordinary diversity in vertebrate populations, is still not fully understood, despite being debated for over four decades (Takahata and Nei, 1990; Spurgin and Richardson, 2010). The main theories comprise pathogen-mediated processes (Spurgin and Richardson, 2010), sexual selection (Milinski, 2006) and maternal-foetal interactions (Hamilton and Hellström, 1978).

Sexual selection or disassortative mating refers to the observation that females produce a disproportionately large fraction of MHC heterozygous offspring by preferentially mating with MHC-dissimilar males, with some species assumed to use body odours as indicators (Apanius et al., 1997; Hughes and Yeager, 1998). This was first observed in the house mouse (Yamazaki et al., 1976), reviewed in Penn (2002) and humans (Wedekind et al., 1995), reviewed in Havlicek and Roberts (2009), and more recently in other vertebrates such as bird species (Bonneaud et al., 2006) and fish species (Consuegra and Garcia de Leaniz, 2008). However, there are also counterexamples, as in studies on wild primates (Huchard et al., 2010) and South Amerindians (Hedrick and Black, 1997), where no such mating preference could be found. Furthermore, it is hard to imagine how MHC-based mate choice would lead to natural selection focused specifically on the PBR (Hughes and Yeager, 1998).

Maternal-foetal interactions by preferential abortion of MHC-homozygous foetuses are another way in which MHC diversity and high levels of heterozygosity could be maintained. However, the empirical evidence for this mechanism is mixed (Hamilton and Hellström, 1978; Hings and Billingham, 1985), and it can only be applied to viviparous animals but not to fish, amphibians, and birds, which have a high MHC diversity as well. Thus, maternal-foetal interactions cannot be seen as a general explanation for MHC diversity (Apanius et al., 1997).

Both these MHC-based reproductive mechanisms not only help maintain genetic diversity within a population, thereby reducing the risk of inbreeding depression (Potts et al., 1994), but also benefit from having MHC-heterozygous offspring, which might be more resistant to pathogens (McClelland et al., 2003). This is because the variability of MHC molecules is correlated with the diversity of the T-lymphocyte receptors, which in turn determine the disease and parasite resistance of an organism and thus can influence the long-term survival probability of populations (Paterson et al., 1998; Sommer, 2005). This now leads to the third and arguably most important mechanism of balancing selection on the MHC, the pathogen-driven (Jeffery and Bangham, 2000) or pathogen-mediated (Spurgin and Richardson, 2010) selection process.

It has long been assumed that pathogens play a central role in evolution (Jeffery and Bangham, 2000), and given the important role the MHC plays in immune recognition, together with the high rate of nonsynonymous substitutions in the PBR which favours amino acid replacements at the PBR positions, it comes as a natural conclusion that pathogen-mediated selection acts as a driving force in maintaining diversity at MHC loci (Doherty and Zinkernagel, 1975; Radwan et al., 2010; Spurgin and Richardson, 2010). However, this selection pressure could act in diverse ways. Three main hypotheses of balancing selection on the MHC mediated by pathogens have been suggested (Spurgin and Richardson, 2010), heterozygote advantage (Doherty and Zinkernagel, 1975), rare allele advantage (Wright and Dobzhansky, 1946; Slade and McCallum, 1992) and selection that varies in time and space (Hill et al., 1991).

The idea of heterozygote advantage (also called overdominance) emerged when Doherty and Zinkernagel (1975) argued that individuals heterozygous at MHC loci can more effectively protect themselves from infections by responding to a broader range of pathogens, as specific MHC-peptide combinations are required to initiate an immune response (Meyer and Thomson, 2001). The heterozygote advantage hypothesis therefore refers to the assumption that heterozygotes are better protected against diseases by responding to a greater range of pathogen peptides than homozygotes (Spurgin and Richardson, 2010) because MHC-mediated resistance is generally dominant (Wakeland et al., 1990; McClelland et al., 2003). This sometimes leads to confusion, since these peptides may be being viewed as belonging to different pathogens (Doherty and Zinkernagel, 1975) or to the same pathogen (O'Connor et al., 2010). The first view is sometimes linked to “dominant selection”, meaning that for a particular pathogen, a heterozygous individual is as fit as the fittest homozygous individual that carries any of the two alleles of the heterozygote genotype, but not fitter (Spurgin and Richardson, 2010). This type of selection would not be able to maintain MHC diversity (Hughes and Nei, 1988). The second viewpoint, however, commonly referred to as “overdominant selection”, where the heterozygote is fitter than both homozygotes, can contribute to MHC diversity (Takahata and Nei, 1990; McClelland et al., 2003), as it increases the mean number of alleles at a locus and the expected heterozygosity (Kimura and Crow, 1964; Maruyama and Nei, 1981), both features of MHC polymorphism.

The hypothesis that overdominant selection maintains the polymorphism of the MHC was initially not widely accepted by population geneticists. Their scepticism was based on theoretical models that did not take the role of mutation in incorporating new alleles into account (Maruyama and Nei, 1981; Hughes and Yeager, 1998), and resulted in no or only a weak polymorphism for most combinations of selection coefficients (Lewontin et al., 1978), with the exception of overdominance that is symmetric in the fitness of heterozygotes.

Negative frequency-dependent selection (Wright and Dobzhansky, 1946; Slade and McCallum, 1992), also referred to as “rare allele advantage” (Wakeland et al., 1990; Meyer-Lucht

and Sommer, 2005) or “minority advantage” (Takahata and Nei, 1990), is another hypothesis that has the potential to maintain MHC diversity. It can be envisaged as a co-evolutionary arms race between hosts and pathogens (Spurgin and Richardson, 2010), leading to a dynamic change in allele frequencies as opposed to a stable equilibrium situation. Pathogens that can evade the immune response of the most common alleles within a population have a selective advantage, since they can thrive within a large share of the host population (Apanius et al., 1997). An allele that is rare or novel in a population would then have high fitness, as no or just few pathogens are able to infect genotypes carrying this allele (Meyer and Thomson, 2001). This leads to an increase in frequency of the formerly rare allele. However, as soon as this allele becomes frequent enough, mutations in the pathogen strain can produce variants that are adapted to this allele. This ultimately leads to a cyclical process, with allele frequencies fluctuating temporally as pathogens adapt to them (Meyer and Thomson, 2001). This form of balancing selection has received a great deal of attention in the literature of theoretical population genetics. For some models of negative frequency-dependent selection it was demonstrated that they are theoretically capable of maintaining a high level of polymorphism, just as some models of heterozygote advantage are (Takahata and Nei, 1990).

It is worth noting that rare alleles have an intrinsic advantage under the heterozygote advantage hypothesis as well, since common alleles will have higher occurrences within homozygotes, which are less fit according to the heterozygote advantage hypothesis. This results in a lower marginal fitness for these alleles. Infrequent alleles however rarely occur in homozygotes, hence their marginal fitness is higher owing to their rareness. This fact makes it particularly hard to distinguish between negative frequency-dependent selection and heterozygote advantage.

A third mechanism that is potentially capable of maintaining MHC diversity is selection that varies in time and space (Hill et al., 1991; Hedrick, 2002). Given that the pathogens a host population is exposed to vary spatio-temporally, it is clear that natural selection on the host population, which in this case is just directional selection, will vary as well (Spurgin and Richardson, 2010). This could result in different MHC alleles being favourable at different points in time and in different areas, hence maintaining diversity across the population.

These three hypotheses are often seen as competing alternatives, hence the scientific debate about which one to favour has been intense for several decades (Lewontin et al., 1978; Spencer and Marks, 1988; Takahata and Nei, 1990; Hughes and Nei, 1992; Slade and McCallum, 1992; De Boer et al., 2004). Surprisingly, even with advanced genetic sequencing methods being available, no consensus has been reached, and the diversity of the MHC still can be seen as one of the major questions in immunogenetics. Identifying the underlying mechanisms of MHC polymorphism and their relative importance would however have major benefits in areas as different as selective breeding (Stear et al., 2005) and conservation

genetics (Sommer, 2005).

Finding out which mechanism of balancing selection operates to what extent within a specific host-parasite system has turned out to be challenging, even when just focusing on parasite-mediated selection mechanisms. It is plausible that in some species all the above mechanisms contribute to MHC polymorphisms (Apanius et al., 1997). Some authors even doubt that identifying the relative importance of these mechanisms within a system is possible (Spurgin and Richardson, 2010). Nevertheless, various attempts have been made, but the results produced an unclear overall picture.

McClelland et al. (2003) performed a coinfection experiment with mice that suggested heterozygote advantage in its overdominant form, whereas Penn et al. (2002) conducted a similar experiment, but only found evidence for dominant-type heterozygote advantage. In a study on red junglefowl, Worley et al. (2010) came to a similar conclusion. However, Ilmonen et al. (2007), who investigated deliberate *Salmonella* infections of a cross of laboratory mice with wild ones, even found that resistance was mostly recessive rather than dominant, and heterozygosity reduced fitness. Heterozygote advantage has been found in two recent studies on Salmonids (Evans and Neff, 2009; Kekäläinen et al., 2009), while another study in a wild salmon population only found allele effects, but did not find any evidence for heterozygote advantage (Dionne et al., 2009) – the authors in fact suggested one of the other pathogen-driven balancing selection mechanisms was occurring. By examining free-ranging mice with nematode infections, Meyer-Lucht and Sommer (2005) also found allele effects without any indication of heterozygote advantage, and suggested that rare allele advantage is the most likely explanation for maintenance of the MHC polymorphism.

On the contrary, Stear et al. (2005) found evidence for heterozygote advantage in the sheep-nematode system, while not excluding the existence of other balancing selection mechanisms. Other work on non-human primates (Sauermaun et al., 2001; O'Connor et al., 2010) and humans (Thursz et al., 1997) also suggests the existence of heterozygote advantage. Finally, studies on mouse lemurs (Schad et al., 2005), primates in their own right, came to the conclusion that the mechanism operating on the MHC has to be frequency-dependent selection rather than heterozygote advantage. However, distinguishing between these two hypotheses proves hard for low sample sizes, as certain (and in particular homozygous) genotypes might be absent from the sample, or only occur at very low frequencies.

1.2 Modelling the selective forces that maintain MHC polymorphism

Modelling can be of help when the empirical evidence is inconclusive. The aim of modelling MHC polymorphism is to show under which conditions - given certain assumptions about alleles and genotypes - this polymorphism may be lost or maintained, and if it is maintained, to what extent.

Earlier models mostly focused on simpler equilibrium situations, as computing power was in short supply. Both Gillespie (1977) and Lewontin et al. (1978) examined situations where a polymorphism of n alleles is stable. The results showed that in a stable environment, the fraction of stable n -allele systems of the total set of all possible fitnesses for the $\frac{n(n+1)}{2}$ genotypes at a locus is vanishingly small for larger n unless the ratio of heterozygote advantage to standard deviation of fitness among heterozygotes is large, in which case the system is close to symmetric overdominance and all alleles have similar frequencies close to $\frac{1}{n}$. If the environment however changes temporally and spatially, then a large number of alleles can be readily maintained, as Gillespie (1977) was able to show analytically for a specific overdominance model. Hoekstra et al. (1985) however pointed out that popular equilibrium models that incorporate heterozygote advantage as well as spatial and temporal variation may not be very robust in terms of maintaining their stability when parameters are varied. Over ten years later, when the variation in space and time was already established as a possible maintenance mechanism for the MHC, Hedrick (2002) developed a model that demonstrated how spatio-temporal fluctuations in pathogen resistance, without any frequency dependent component or intrinsic heterozygote advantage, can maintain a stable polymorphism. This model incorporated population genetic results, while Monte Carlo simulations generated distributions of allele frequencies.

Heterozygote advantage, the first mechanism suggested to maintain the MHC polymorphism, has been modelled extensively. Using stochastic differential equations, Maruyama and Nei (1981) found that overdominant selection is a powerful mechanism in increasing the mean heterozygosity compared with a neutral model, which let them conclude that overdominance can explain a higher degree of MHC diversity. Spencer and Marks (1988) discussed the maintenance of MHC polymorphism at a single MHC locus when allowing for mutations, using Monte Carlo simulations. These mutations resulted in alleles with random fitnesses. Their model showed that selection is capable of maintaining much more alleles in a population than predicted by Lewontin et al. (1978), as unstable fitness matrices, and thus alleles that confer low fitness, are quickly eliminated by natural selection. Using deterministic simulations, Marks and Spencer (1991) were able to refine their model while highlighting the significance of how the mutation process is modelled, and in a further refinement (Spencer

and Marks, 1992), they could also reach a higher degree of MHC polymorphism with their simulations, bringing MHC diversity predicted by the model closer to observed values.

Using Monte Carlo simulations that included mutations, Takahata and Nei (1990) investigated alternative hypotheses for mechanisms maintaining MHC polymorphism. They developed a computer simulation model that not only investigated the conditions for a stable polymorphism, but also reconstructed genealogic relationships of alleles, as they considered the latter as an important property of any model examining MHC diversity. The results of these simulations suggested that in the symmetric overdominance model, which assumed that all heterozygotes have equal fitness (they were assigned a fitness value of 1), while all homozygotes have lower fitness, a polymorphism at least of the order of values observed in nature can easily be maintained, again in stark contrast to the conclusions of Lewontin et al. (1978) for the equilibrium case without mutation. Nevertheless, the degree of polymorphism decreased with increasing asymmetry of the fitness values of homozygotes, just as in the equilibrium analysis. In simulations of Takahata and Nei (1990), the frequency-dependent selection model delivered essentially the same results as the overdominance selection models, so neither model could be discriminated against on mathematical grounds; the authors favoured the overdominance model for biological reasons.

To assess the combined effect of intra-locus recombination and balancing selection, Satta (1997) performed computer simulations of two models of balancing selection, a symmetric overdominance model and an asymmetric model that was based on the number of codon differences between the two alleles on the examined locus. The model repeated the processes of mutation, sampling, and intra-locus recombination in each time step. Comparing the simulation results to observed data, Satta (1997) concluded that symmetric overdominance is more appropriate than the asymmetric, sequence-based model, but that the assumed locations of recombination break points were of great importance.

More recently, De Boer et al. (2004) constructed an allele-fitness model by assigning fitness values to alleles, which represented the fraction of pathogens to which a particular allele can provide protection. Assuming an additive allele effect on the genotype with average fitness overlap of the alleles and using an earlier result of conditions for the persistence of alleles under general conditions (Lieberman, 1991), they conceived an asymmetric overdominance model where allelic persistence at equilibrium could be determined using a simple formula. Furthermore, De Boer et al. (2004) performed stochastic computer simulations for a specific case of their model to demonstrate that the same result holds for a finite population and with mutations. Both approaches suggested that heterozygote advantage alone is insufficient to maintain the amount of MHC diversity observed. In a companion paper, Borghans et al. (2004) again performed simulations, which coded peptides as binary bit strings, and showed that a frequency-dependent model of host-pathogen coevolution is able to maintain the observed levels of MHC polymorphism. This model had been adapted by Ejsmond et al.

(2010), who found that observed allele frequencies are often indistinguishable from neutrality expectations under negative frequency-dependent selection.

The models discussed so far explored the classic hypotheses of pathogen-mediated selection. There are, however, alternative approaches worth mentioning. While classic heterozygote advantage assumes a higher fitness for heterozygotes, as heterozygotes have different alleles and can therefore recognise more pathogens, Stoffels and Spencer (2008) speculated that MHC genotypes might have a higher fitness if a pathogen belonged to the portion where the recognition sets of both alleles overlapped. While no experimental evidence of such an “intersection advantage” exists, they concluded nevertheless that if it existed, heterozygote advantage might not be sufficient to explain the high levels of polymorphism observed at MHC genes. Finally, van Oosterhout (2009) introduced a new model of MHC evolution named “Associative Balancing Complex” (ABC) evolution, that, according to the author, resolved deficiencies of traditional models of balancing selection (as overdominance and negative frequency dependent selection). This model assumed that if a certain degree of diversity already existed at the MHC, then recessive deleterious mutations could accumulate as a “sheltered load” nearby MHC genes. If a haplotype now became very common so that it regularly formed homozygotes, then purifying selection would reduce its frequency, thus maintaining a balanced MHC polymorphism even without strong selection by parasites.

Decades of research on the causes of the polymorphism of genes encoding the MHC have not been able to unequivocally ascertain the relative importance of different diversity maintenance mechanisms within a particular host-parasite system, or across systems. Even the large body of data obtained from non-model species in natural populations was not sufficient to identify the processes that underpin selection at MHC genes (Piertney and Oliver, 2006). Many studies suffer from the fact that they have not fully considered the alternative explanations of possible balancing selection mechanisms (Spurgin and Richardson, 2010). An additional problem when testing different hypotheses of balancing selection is the large sample size needed for a conventional population study, as in many species most individuals are heterozygous at most MHC loci, so a large enough sample of homozygotes is hard to obtain (Hughes and Yeager, 1998).

This is where modelling can make important contributions. For example, models can be used to establish conditions for the maintenance of diversity, which might be easier to check in a field study than the polymorphism itself. However, modelling approaches undertaken so far have delivered conflicting results, depending on the assumptions and emphasis of the model. Therefore a more generic approach might be necessary to improve our understanding of the maintenance mechanisms of this peculiar genetic region that exhibits an unrivalled diversity.

1.3 Structure and function of the MHC

Structural information on MHC molecules is important to evolutionary studies, as it offers predictions about how selection acts on functionally distinct regions of a MHC molecule. Selection favouring increased scope of pathogen presentation, which results from a high rate of nonsynonymous substitutions, is concentrated at sites involved in peptide presentation (the PBR), whereas changes at sites that are critical for tertiary structure impair the molecule's function, so selection eliminates mutants and keeps such sites largely unmodified (Meyer and Thomson, 2001). This makes it possible to partition codons based on their selective pressures.

The major histocompatibility complex (MHC) of vertebrates is a multigene family of closely linked genes acting at the interface between the immune system and infectious diseases (Bernatchez and Landry, 2003). MHC loci encode receptors (glycoproteins) on the surface of various immune and nonimmune cells (Bernatchez and Landry, 2003), making the MHC the most important genetic component of the mammalian immune system (Klein, 1986), the “center of the immune universe” (Edwards and Hedrick, 1998).

The primary role of the MHC is the recognition of foreign proteins originating from pathogens, the presentation of short fragments (“peptides”) of these proteins to certain immune cells and thus the initiation of an immune response (Klein, 1986). More specifically, foreign proteins, which might have entered cells by infection or phagocytosis, are broken down into peptides and loaded onto specific MHC molecules (Pierny and Oliver, 2006). Some of these complexes of MHC and peptides are subsequently transported to the cell surface and presented to the circulating immune cell (specifically T-cell) population. The immune response is triggered if the MHC-peptide complex is recognized as foreign (Meyer and Thomson, 2001) and T-cells bind to the presented peptide (Pierny and Oliver, 2006). The immune response is thus MHC restricted in the sense that T-cell recognition of infected cells requires a combined signal from both MHC molecules and pathogen peptides, and this restriction makes histocompatibility molecules key players in the adaptive immune response (Meyer and Thomson, 2001).

Because of the importance of the MHC molecules and their role in the immune response, an extensive amount of information is available about structure and function of MHC molecules (Trowsdale, 1993). The extracellular components of MHC molecules are composed of two main parts, a “stalk” that anchors the molecule in the surface of cells, and a groove-shaped receptor, the PBR (Edwards and Hedrick, 1998), which holds peptides in the groove by noncovalent interactions with amino acid residues (Klein et al., 1993). It is the PBR that is responsible for the recognition of antigens, and a match between PBR, antigenic peptide and T-cell receptor is required to initiate the complex cascade of immune responses (Pierny and Oliver, 2006). Although PBRs are specific to some extent, a single MHC molecule can still

bind to a range of different self- and non-self peptides (Altuvia and Margalit, 2004), given that these have specific amino acids at certain anchor sites that fit exactly to the binding groove of the MHC molecule (Spurgin and Richardson, 2010). The amino acids lining the groove of the PBR determine the antigen-binding specifics of the MHC molecule, i.e. the association between MHC molecule and peptides that it can bind (Hedrick, 1994).

MHC molecules are members of the large immunoglobulin superfamily of immunologically active molecules (Edwards and Hedrick, 1998; Piertney and Oliver, 2006), and can be further subdivided into two major subfamilies, called class I and class II (Hughes and Yeager, 1998). The three-dimensional structure of both human MHC class I (Bjorkman et al., 1987) and class II (Brown et al., 1993) molecules are known. This knowledge gives extraordinary insight into the main function of these molecules (Piertney and Oliver, 2006), which is the recognition and presentation of short peptides, 8 to 11 amino acids in length for class I molecules (Meyer and Thomson, 2001) and 11 to 25 for class II (Hedrick, 1994; Rammensee et al., 1995; Meyer and Thomson, 2001). The reason for the much lower variability of class I molecules is that the ends of the peptide are tucked down into the peptide-binding groove, which limits the peptide's length. In class II, the peptide's ends are free, and therefore longer peptides are possible (Rammensee et al., 1995; Hughes and Yeager, 1998).

MHC class I genes are expressed on the surface of most nucleated somatic cells (Bernatchez and Landry, 2003; Piertney and Oliver, 2006). They play an essential role in the immune defence against *intracellular* pathogens by binding peptides mainly derived from viral proteins or malignant cells (Hedrick, 1994; Bernatchez and Landry, 2003; Sommer, 2005) and presenting them to cytotoxic T-cells (Piertney and Oliver, 2006). Class I MHC molecules are heterodimers, consisting of three extracellular domains (α_1 , α_2 , α_3) that form the "class I heavy chain" and β_2 microglobulin (Jeffery and Bangham, 2000). Each of the extracellular domains is encoded by a different exon from a single gene (Piertney and Oliver, 2006). The PBR of class I molecules consists of a groove formed by the α_1 and α_2 domains, sitting on top of a β -pleated sheet (Bjorkman et al., 1987; Edwards and Hedrick, 1998; Jeffery and Bangham, 2000).

MHC class II genes are predominantly expressed on antigen-presenting cells of the immune system, such as B cells, macrophages, lymphocytes and dendritic cells (Bernatchez and Landry, 2003; Piertney and Oliver, 2006). Class II molecules present processed exogenous antigens derived from *extracellular* parasites and pathogens (as, for example, bacteria and nematodes) to T-helper cells (Bernatchez and Landry, 2003; Piertney and Oliver, 2006), which in turn release cytokines that trigger an appropriate immune response, as the production of antibodies and stimulation of macrophages (Hughes and Yeager, 1998; Meyer and Thomson, 2001). Class II molecules are heterodimers consisting of an α chain and a β chain (Hughes and Yeager, 1998). These two chains together form a peptide-binding groove at the top of the molecule, which is similar in structure to that of the class I molecule (Brown et al.,

1993). The class II region of placental mammals is divided into subregions, each of which contains an α chain gene and one or more β chain genes. In humans, these subregions are designated DR, DP, and DQ (Hughes and Yeager, 1998). Specific sites within both the α chain and the β chain form the class II PBR (Pirotney and Oliver, 2006).

Sequence variation among MHC alleles of both class I and II genes that encode parts of the PBR is assumed to enable recognition of various pathogens (Potts and Wakeland, 1990). The peptides that a given MHC molecule binds share amino acids at a number of positions (the anchor residues). These positions determine the specificity of binding by interacting with a subset of the amino acids which make up the PBR of the MHC molecule. A single MHC molecule can thus bind many peptides, if these peptides only have the correct amino acids at the anchor positions (Meyer and Thomson, 2001). Nonetheless, because different allelic products have different anchor motifs, a heterozygote should have an advantage over a homozygote in terms of immune surveillance.

1.4 Measuring diversity

Discussing selective forces that maintain diversity at the MHC raises the question how diversity should be measured. No precise definition has been presented so far, even though a clear concept of diversity is highly desirable when comparing predictions of different models, or model predictions to observed data. One of the most intuitive notions is perhaps the approximation of MHC diversity by the number of alleles that are maintained, or that exist at a particular point in time. Many authors, in the absence of reliable data on allele frequencies, have used the number of alleles as the main or even sole criterion for determining the validity of a particular hypothesis.

Gillespie (1977) and Lewontin et al. (1978) explored conditions for a stable polymorphism of n alleles, while Hoekstra et al. (1985) determined under which conditions such an n -allele polymorphism is robust, thus neglecting the frequencies of the persisting alleles. Spencer and Marks (1988), while mainly concerned about the number of alleles that could be maintained in their Monte Carlo simulations, mentioned frequency distributions of alleles, while Maruyama and Nei (1981) additionally compared the (mean) heterozygosity, and Hedrick (2002) as well as Stoffels and Spencer (2008) compared observed and expected heterozygosity values to theoretical ones. Marks and Spencer (1991) reported the number of alleles from other studies, but in their own work appreciated the importance of allele frequency distributions, and how even or uneven they were, albeit without measuring the degree of evenness. In their companion paper, Spencer and Marks (1992) introduced an index that measured a distance of allele frequency vectors from a vector of uniform allele frequencies, therefore quantifying evenness. Heterozygosity and the number of alleles are again the key

witnesses when Takahata and Nei (1990) compared allele frequency distributions of simulations to observed data, while De Boer et al. (2004) compared distributions of frequencies and fitness contributions of alleles without using a measure to quantify them. The unevenness of allele frequency distributions was finally used by van Oosterhout (2009) as supporting evidence for the validity of his model of “ABC evolution” (see section 1.2), but he did not show how well the degree of unevenness matched that of observed data.

All these examples show that the debate about selective forces maintaining MHC diversity does not fully acknowledge the need to further specify diversity. “Diversity” is rather used freely, often with different meanings. I will now introduce a powerful method that gives a more general picture of diversity than single indices (as the number of alleles or the expected heterozygosity) can. This method produces a diversity profile, i.e. a continuum of diversity measures, by varying a “viewpoint parameter” from zero to infinity. This results in a “fingerprint” of the gene pool of a population, and additionally encompasses traditional diversity indices such as the number of alleles, expected homo- and heterozygosity or the frequency of the most common allele. It can further be naturally expanded to include similarity data, which is a great asset when such data can be sensibly defined for MHC alleles. Diversity profiles will be applied to allele frequency and similarity data in the subsequent chapters.

For calculating diversity profiles, I used an approach described by Leinster and Cobbold (2012) that calculates diversity measures ${}^qD^Z$, where q is the viewpoint (or sensitivity) parameter that is varied along the real axis, and Z is a matrix that contains the measure of similarity for every pair of alleles. When similarity between alleles is not taken into account, which is subsequently called “naive” diversity, Z equals the identity matrix I .

This method comes with a number of advantages. First, it is more informative than just a single statistic, and replaces a number of popular indices of allelic diversity by a single formula. Second, it behaves intuitively as it uses effective numbers, so that percentage changes and ratio comparisons of diversity are meaningful. Third, it allows for useful comparisons of diversity both with or without taking similarity between MHC alleles into account, and fourth, it is highly versatile, since it allows similarity to be measured in different ways according to the current objective and the availability of particular types of similarity data.

The diversity profile ${}^qD^Z(\mathbf{p})$, $q \in [0, \infty]$ takes two inputs:

1. Frequency data of alleles (the frequencies of all S alleles in the gene pool sum to one, $\sum_{i=1}^S p_i = 1$)
2. Similarity data: for each pair of alleles, a number between zero and one that specifies how “similar” (in terms of the quantity of interest) the alleles are.

The viewpoint parameter q ($0 \leq q \leq \infty$) indicates how significant the abundance of alleles is: for $q = 0$, rare alleles are treated equally to common ones, i.e. frequencies of alleles

are entirely disregarded, while for $q = \infty$, diversity depends only on the most frequent allele; rare alleles are ignored altogether. Given frequency and possibly similarity data, the diversity of order q is calculated for every q , and the result is plotted against q . This graph is the diversity profile of the set of alleles in the current gene pool; it contains information that can be read off at a glance.

Below I give a brief overview of the calculation method for the diversity profile ${}^qD^Z$ from Leinster and Cobbold (2012), where a more detailed description can be found.

For a gene pool of S alleles, the allele frequencies, denoted by p_1, \dots, p_S , are combined to a frequency vector $\mathbf{p} = (p_1, \dots, p_S)$. Similarities between alleles are encoded in an $S \times S$ matrix $\mathbf{Z} = (Z_{ij})$, with Z_{ij} measuring the similarity between the i th and j th allele, and $0 \leq Z_{ij} \leq 1$, with 0 indicating entirely dissimilar alleles and 1 indicating identical alleles (hence $Z_{ii} = 1$). There is one diversity measure for each value of the viewpoint parameter q , which controls the relative emphasis on common and rare alleles.

For $q \neq 1, \infty$, the *diversity of order q* of the community is defined as

$${}^qD^Z(\mathbf{p}) = \left(\sum_{i|p_i>0} p_i (\mathbf{Z}\mathbf{p})_i^{q-1} \right)^{\frac{1}{1-q}} \quad (1.1)$$

where

$$(\mathbf{Z}\mathbf{p})_i = \sum_{j=1}^S Z_{ij} p_j$$

The cases $q = 1$ and $q = \infty$ have been excluded because equation 1.1 for ${}^qD^Z(\mathbf{p})$ does not make sense there. It does, however, converge to a limit as $q \rightarrow 1$ or $q \rightarrow \infty$, therefore ${}^1D^Z(\mathbf{p})$ and ${}^\infty D^Z(\mathbf{p})$ are defined as those limits:

$${}^1D^Z(\mathbf{p}) = \frac{1}{(\mathbf{Z}\mathbf{p})_1^{p_1} (\mathbf{Z}\mathbf{p})_2^{p_2} \dots (\mathbf{Z}\mathbf{p})_S^{p_S}} \quad (1.2)$$

$${}^\infty D^Z(\mathbf{p}) = \frac{1}{\max(\mathbf{Z}\mathbf{p})_i} \quad (1.3)$$

The diversity profile of the set of alleles examined is given by the graph of ${}^qD^Z$ against q . The left-hand end of a diversity profile (when q is small) is affected almost as much by rare alleles as common ones, whereas the right-hand tail (when q is large) gives information about common alleles, but largely disregards rare ones (Hill, 1973). As the viewpoint parameter q grows, the perceived diversity ${}^qD^Z(\mathbf{p})$ drops, and therefore the diversity profile is always a decreasing curve. The values of ${}^0D^Z$, ${}^1D^Z$, ${}^2D^Z$ and ${}^\infty D^Z$ usually give a good indication of the shape of this curve.

Chapter 2

The Maintenance of MHC Alleles at Equilibrium – A General Framework

2.1 Introduction

Using classical single locus multi-allele viability models (Lewontin et al., 1978; Karlin and Lessard, 1984), a general framework for classifying models of heterozygote advantage is introduced. It is demonstrated how existing models of heterozygote advantage, symmetric overdominance (Robertson, 1962; Takahata and Nei, 1990; Satta, 1997) and asymmetric overdominance (Satta, 1997; De Boer et al., 2004), represent special cases of this framework.

I review methods to calculate the number of persisting alleles and their equilibrium frequencies, and discuss the stability of this equilibrium, before applying these results to alleles of the MHC. Using the divergent allele advantage hypothesis proposed by Wakeland et al. (1990), a novel model is presented, subsequently called the divergent allele advantage model; it is shown how this novel model fits into the general framework discussed before, and how it relates to traditional models of heterozygote advantage.

All these considerations are based on the assumption that the system (which is in general a set of alleles that occurs in a population at a specific locus) is at equilibrium, and that there is no spatial or temporal variation in the environment. The equilibrium assumption will only be dropped in chapters 4, 5 and 6, where new alleles can be created by mutation and lost by genetic drift, and therefore a real equilibrium can generally never be reached. Chapter 6 additionally allows for spatial variation and migration by exploring metapopulation dynamics.

2.2 Materials and Methods

2.2.1 General description of a single locus model

I consider an infinite population with discrete, non-overlapping generations and random mating, and examine a single autosomal locus with alleles A_i ($i = 1, 2, \dots, k$), which occur at relative frequencies p_i ($\sum_{i=1}^k p_i = 1$) within the population. Assuming the population is in Hardy-Weinberg equilibrium, then the genotype of an individual with alleles A_i and A_j ($i \leq j$), $A_i A_j$, occurs at the frequency $2p_i p_j$ for $i < j$ and p_i^2 for $i = j$. Since there is little evidence of segregation distortion at the MHC genes, I define the fitness of an individual $A_i A_j$, f_{ij} , as the offspring contribution this individual can make to the next generation. As my main focus is on the pathogen-mediated processes in a host-pathogen system, I consider the fitness of a genotype as the effectiveness with which the host's immune system responds to the challenge of the pathogens it is exposed to. A fitness value of 0 corresponds to genotypes that are not viable, whereas a fitness value of 1 corresponds to genotypes that are fully protected against all pathogens. The classical single locus multi-allele viability model of population genetic theory (Crow and Kimura, 1970) leads to an update formula for calculating the proportions of all alleles in the system in the next generation, given that the proportions in the current generation are known. In matrix formulation, this can be written as (Lewontin et al., 1978; Karlin and Lessard, 1984):

$$\mathbf{p}_{t+1} = \frac{\mathbf{p}_t \circ \mathbf{F} \mathbf{p}_t}{\bar{w}(\mathbf{p}_t)}, \quad \bar{w}(\mathbf{p}_t) = \sum_{i=1}^k \sum_{j=1}^k f_{ij}(\mathbf{p}_t)_i (\mathbf{p}_t)_j \quad (2.1)$$

Here, \mathbf{p}_t and \mathbf{p}_{t+1} are the proportions of all alleles in the system at times t and $t + 1$, $\mathbf{F} = (f_{ij})_{1 \leq i, j \leq k}$ is the genotype fitness matrix of all genotypes made up of alleles from the set A_1, \dots, A_k , \circ denotes component-wise multiplication of two vectors, and $\bar{w}(\mathbf{p}_t)$ is the population fitness at time t . In other words, the frequency of A_i at time $t + 1$ can be calculated by multiplying the frequency of this allele at time t by the relative marginal fitness of the same allele:

$$(\mathbf{p}_{t+1})_i = (\mathbf{p}_t)_i \cdot \frac{(w_m(\mathbf{p}_t))_i}{\bar{w}(\mathbf{p}_t)} \quad (i = 1, 2, \dots, k) \quad (2.2)$$

The marginal fitness of an allele is defined as the average fitness of the genotypes in which it is present, weighted by the proportion of each genotype in the population:

$$(w_m(\mathbf{p}_t))_i = (\mathbf{F} \mathbf{p}_t)_i = \sum_{j=1}^k f_{ij}(\mathbf{p}_t)_j \quad (i = 1, 2, \dots, k) \quad (2.3)$$

The marginal fitness of an allele is hence closely related to the quantitative genetics concept of the additive value (in terms of fitness), as the marginal fitness is equal to the mean fitness of the genotypes this allele produces.

The population fitness $\bar{w}(\mathbf{p}_t)$ can also be expressed in terms of the marginal fitness values of the alleles:

$$\bar{w}(\mathbf{p}_t) = \sum_{j=1}^k (\mathbf{p}_t)_j \cdot (\mathbf{w}_m(\mathbf{p}_t))_j \quad (2.4)$$

Equation 2.2 describes a discrete time dynamical system. To find the allele frequencies at equilibrium, the fixed points of this system have to be found.

If a vector \mathbf{x} with $\sum_{j=1}^k x_j \neq 0$ exists that solves the system of linear equations

$$\sum_{j=1}^k f_{ij} x_j = 1 \quad (i = 1, 2, \dots, k) \quad (2.5)$$

then the k -allele system is called internally stable if every component of the vector \mathbf{x} is positive. In this case the vector of equilibrium frequencies \mathbf{p} can be calculated by normalising the components of \mathbf{x} :

$$p_i = \frac{x_i}{\sum_{j=1}^k x_j} \quad (i = 1, 2, \dots, k) \quad (2.6)$$

This vector \mathbf{p} contains the equilibrium frequencies, as 2.3 and 2.5 imply that

$$(\mathbf{w}_m(\mathbf{p}_t))_i = \sum_{j=1}^k f_{ij} \frac{x_j}{\sum_{l=1}^k x_l} = \left(\sum_{l=1}^k x_l \right)^{-1} \cdot \sum_{j=1}^k f_{ij} x_j = \left(\sum_{l=1}^k x_l \right)^{-1} \quad (i = 1, 2, \dots, k) \quad (2.7)$$

and, according to 2.4, the population fitness is:

$$\bar{w}(\mathbf{p}_t) = \left(\sum_{l=1}^k x_l \right)^{-1} \cdot \sum_{j=1}^k (\mathbf{p}_t)_j = \left(\sum_{l=1}^k x_l \right)^{-1} \quad (2.8)$$

Equation 2.2 together with 2.7 and 2.8 demonstrates that the allele frequencies are indeed at equilibrium:

$$(\mathbf{p}_{t+1})_i = (\mathbf{p}_t)_i \cdot \frac{(\mathbf{w}_m(\mathbf{p}_t))_i}{\bar{w}(\mathbf{p}_t)} = (\mathbf{p}_t)_i \cdot \frac{\sum_{l=1}^k x_l}{\sum_{l=1}^k x_l} = (\mathbf{p}_t)_i \quad (2.9)$$

More details regarding the limitations of this method, particularly with respect to cases where F or a submatrix of it is singular, can be found in Karlin and Lessard (1984) and Liberman (1991). An alternative method of determining the allele frequencies at equilibrium, based on Fisher's fundamental theorem of natural selection, will be discussed in chapter 3.

2.2.2 The genotype fitness matrix

The previous section demonstrated that the system dynamics are ultimately characterised by the genotype fitness matrix \mathbf{F} . This means that the properties of this single locus multi-allele system for different models of overdominance (see section 2.3), such as stable equilibria or allele persistence, depend on the relative magnitude of the elements f_{ij} of \mathbf{F} .

To characterise the genotype fitness matrix \mathbf{F} in terms of the fitness contributions of individual alleles, I define the intrinsic merit w_i of an allele A_i as the fitness of a homozygote that contains two copies of this allele at the considered locus, i.e. $w_i = f_{ii}$. This is tantamount to equating the intrinsic merit w_i with the genotypic value of the homozygote A_iA_i if fitness is the measure of interest. The intrinsic merit of an allele A_i is thus independent of the frequency of the allele A_i , in contrast to the marginal fitness of this allele or the breeding value of the genotype A_iA_i .

I further assume that the intrinsic merits of all different alleles are ordered ($w_1 \geq w_2 \geq \dots \geq w_k$). Allowing for different intrinsic merits for different alleles is based on the observation that alleles are capable of conferring protection from and susceptibility to infectious diseases (Meyer and Thomson, 2001).

As the genotype fitnesses of heterozygotes $f_{ij} (i \neq j)$ depend on the alleles and the overdominance model used, it is necessary to define how two alleles on the considered locus translate to a genotype fitness (or genotypic value of the heterozygote A_iA_j in terms of fitness); this is explained in section 2.3.

The genotype fitness matrix \mathbf{F} is always symmetric, as I assume Mendelian segregation, i.e. (Weissing and van Boven, 2001):

$$f_{ji} = f_{ij} \quad (\forall i, j) \quad (2.10)$$

Thus, in its most general form, the genotype fitness matrix \mathbf{F} can be written as

$$\mathbf{F} = \begin{pmatrix} w_1 & f_{12} & \dots & f_{1k} \\ f_{12} & w_2 & \dots & f_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ f_{1k} & f_{2k} & \dots & w_k \end{pmatrix} \quad (2.11)$$

A fully polymorphic equilibrium (k extant alleles) leads to a positive solution \mathbf{p}^* of the equation systems 2.5 and 2.6 (Weissing and van Boven, 2001). The vector \mathbf{p}^* represents the “optimal frequencies”, i.e. the allele frequencies which correspond to a strict local fitness maximum of the population. Hence, if the stable equilibrium with frequencies \mathbf{p}^* is unique, then $\bar{w}(\mathbf{p}^*)$ is a global maximum.

2.3 Applications to selected overdominance models

2.3.1 Symmetric overdominance

In a symmetric overdominance (SOD) model, all heterozygotes have the same fitness value of 1 ($f_{ij} = 1$ for $i \neq j$, Robertson (1962)), i.e. heterozygotes are fully protected against every pathogen according to the definition of fitness of section 2.2.1. If $w_1 < 1$, then this model represents overdominance: the heterozygote is always fitter than both homozygotes. This leads to the following genotype fitness matrix \mathbf{F} :

$$\mathbf{F} = \begin{pmatrix} w_1 & 1 & \dots & 1 \\ 1 & w_2 & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \dots & w_k \end{pmatrix} \quad (2.12)$$

Here, the k -allele polymorphism $\{A_1, \dots, A_k\}$ is always maintained, independent of the intrinsic merits of the alleles w_i (Marks and Spencer, 1991). What is more, new alleles are always able to invade the polymorphism of k alleles, and they do not displace any of the k original alleles (De Boer et al., 2004). Hence, such a symmetric overdominance model accumulates alleles, and loss of alleles is only possible if the population size is finite and stochastic extinction is allowed. A special case of the overdominance model described here is the case where all alleles have the same intrinsic merit ($w_i = w_j < 1 \quad \forall i, j$), i.e. not only do all heterozygotes have equal fitness, but all homozygotes have equal fitness as well, and the fitness of the homozygotes is lower than that of the heterozygotes (Meyer and Thomson, 2001). I subsequently call such a model a *fully* symmetric overdominance (FSOD) model.

2.3.2 Simple asymmetric overdominance

Simple asymmetric overdominance assumes an equal and constant advantage of heterozygotes over homozygotes (reviewed in Hedrick (1999)). The special case that all homozygotes have the same fitness (i.e. all alleles have the same intrinsic merit) and all heterozygotes have the same elevated fitness,

$$0 \leq f_{ii} = w_i = d < 1 \quad \forall i, \quad f_{ij} = f_{ii} + c = d + c \quad (i \neq j, c > 0) \quad (2.13)$$

is just a fully symmetric overdominance model as described above. In general, however, the intrinsic merits of the alleles w_i do not all need to have the same value, in which case the fitness values of the heterozygotes $f_{ij}, i \neq j$ vary accordingly. In the simple asymmetric overdominance (SAsOD) model explored here I allow for variation in the intrinsic

merit of the alleles, but assume that the amount of heterozygote advantage c is constant for all heterozygous genotypes. This has the effect that both the homozygous and heterozygous genotypes have fitness values that differ among themselves, which makes such a model asymmetric for both homo- and heterozygotes. The heterozygote advantage c is added to the average of the intrinsic merits of the two alleles an individual carries according to equation 2.14:

$$f_{ij} = \frac{1}{2} \cdot (w_i + w_j) + c \quad (i \neq j) \quad (2.14)$$

The results in terms of stability of the k -allele system strongly depend on whether the heterozygote fitnesses f_{ij} may exceed 1 or – in line with the initial notion of representing a fraction of pathogen peptides – are bounded by 1, i.e.

$$f_{ij} = \min\left(\frac{1}{2} \cdot (w_i + w_j) + c, 1\right) \quad (i \neq j) \quad (2.15)$$

If fitnesses represent fractions of pathogen peptides, then only the second option is biologically plausible, and therefore I restrict discussion to this case.

2.3.3 Reinforcing asymmetric overdominance

As in the simple asymmetric overdominance model discussed in section 2.3.2, the fitness of a heterozygous individual of genotype $A_i A_j$ ($i \neq j$) also depends on the intrinsic merits of its alleles in the reinforcing asymmetric overdominance (RAsOD) model. In the RAsOD model, however, the fitness advantage gained by heterozygotes is specified by equation 2.16 (De Boer et al., 2004):

$$f_{ij} = w_i + (1 - w_i) \cdot w_j \quad (i \neq j) \quad (2.16)$$

This can be envisaged as a combined protective effect of both alleles, which is reinforcing in that two different alleles with a low intrinsic merit will not combine to give a heterozygous genotype with above average fitness. On the contrary, the heterozygotes' fitnesses are strictly ordered according to the underlying allele intrinsic merits. This model is also, as the previous one, an overdominance model: in case that $0 < w_k \leq \dots \leq w_1 < 1$, then for each off-diagonal element $\max(w_i, w_j) < f_{ij} = w_i + (1 - w_i) \cdot w_j < 1$ ($i \neq j$). The genotype fitness matrix \mathbf{F} of this model can be written as

$$\mathbf{F} = \begin{pmatrix} w_1 & w_1 + (1 - w_1) \cdot w_2 & \dots & w_1 + (1 - w_1) \cdot w_k \\ w_1 + (1 - w_1) \cdot w_2 & w_2 & \dots & w_2 + (1 - w_2) \cdot w_k \\ \vdots & \vdots & \ddots & \vdots \\ w_1 + (1 - w_1) \cdot w_k & w_2 + (1 - w_2) \cdot w_k & \dots & w_k \end{pmatrix} \quad (2.17)$$

This model leads to a situation where the internal stability of the k -allele system only depends on the intrinsic merits of the alleles. As demonstrated by De Boer et al. (2004), the internal stability of the system can simply be derived by calculating a threshold value

$$t = \frac{k-1}{k} \cdot \hat{w} \quad (2.18)$$

where \hat{w} is the harmonic mean of w_1, \dots, w_n . This immediately allows conclusions to be drawn about internal stability of the k -allele system without having to solve a system of linear equations: the k -allele system is only internally stable if the intrinsic merits w_i of all alleles are above the threshold value.

2.3.4 Divergent allele advantage

Divergent allele advantage (DAA) is a concept first presented by Wakeland et al. (1990). It assumes that evolution favours alleles with highly divergent forms of the antigen binding site. This may result in a better protection against pathogens, since dissimilar alleles have less overlap and therefore greater coverage of the space of pathogen epitopes.

The idea of this concept is used to create a model of overdominance that can be characterised in the same framework as the overdominance models discussed previously. This novel model, inspired by the DAA hypothesis, basically works like a binary AND operation between pathogen peptide recognition sites on the two alleles when the fitness of a heterozygous genotype is determined (figure 2.1). The number of epitopes the heterozygous individual shown in figure 2.1 can recognise is the number of epitopes in the union of the set of epitopes that an allele A recognises and the set of epitopes that an allele B on the same locus of the other chromosome recognises. As in the RAsOD model, an overlap of epitopes (the epitopes both alleles recognise) might exist. The key difference between the two models is that overlap is not calculated as an average property of the alleles, but depends on the alleles involved. Hence, this model could be called a non-reinforcing model: alleles of lower intrinsic merit might still combine in a complementary way (with just little or no overlap) to form heterozygotes with high fitness.

Figure 2.1 provides an example. Here, allele A recognises 5 epitopes out of 10 (epitopes recognised are symbolised by black squares) and has therefore an intrinsic merit of $w_A = 0.5$. Allele B recognises only 4 epitopes out of 10 and has hence an intrinsic merit of $w_B = 0.4$. However, only the epitope on position 3 (from left) is recognised by both alleles, all the other positions are either recognised by allele A only (positions 1, 6, 7 and 10), or by allele B only (positions 2, 4 and 9) or by neither allele (position 5 and 8). As the MHC genes are codominantly expressed (Potts and Slev, 1995), it can be assumed that it is sufficient for one of the alleles to recognise an epitope in order to ensure recognition by the genotype. This

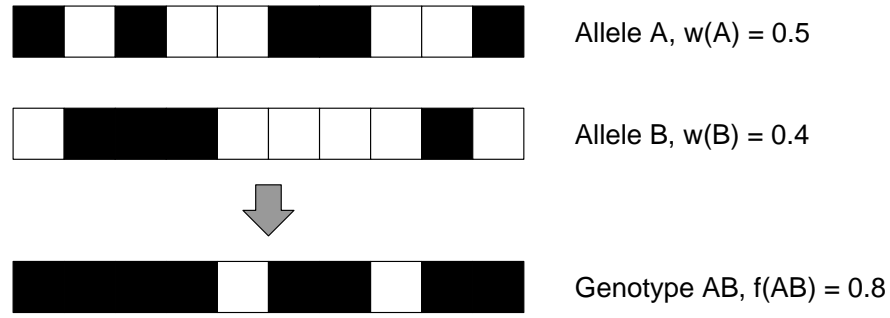


Figure 2.1: Example of the calculation of the genotype fitness from the intrinsic merits of the alleles in the DAA model. Black squares stand for recognised epitopes, whereas epitopes that are not recognised are represented by white squares.

results in a genotype fitness $f_{AB} = 0.8$, as 8 out of 10 epitopes are recognised by either allele A or allele B .

In the genotype fitness matrix framework, the off-diagonal elements of this model will on average be larger compared to the asymmetric overdominance models if strong divergence between the alleles is assumed, even if the intrinsic merits of the alleles are the same in all examined models. Hence, a stabilising effect on the k -allele system can be expected under DAA, which in turn can increase both the number of alleles that can be maintained at equilibrium, and the difference in intrinsic merit between the best and the worst extant allele.

The DAA model can only maintain a higher number of alleles than the RAsOD model if the overlap between the alleles in the system is sufficiently small, which in turn results in a good coverage of the pathogen epitope space. This is demonstrated in the following simple example: Assume a system of three alleles, with intrinsic merits of $w_1 = 0.8$, $w_2 = 0.7$, $w_3 = 0.1$. The RAsOD model predicts that such a system is unstable, since the intrinsic merit of allele A_3 , w_3 , lies below the stability threshold value of $\frac{k-1}{k} \cdot \hat{w} \approx 0.158$. In the DAA model, however, the pathogen recognition sites of both alleles need to be known in order to calculate the fitnesses of the genotypes, f_{ij} , which in turn can be used to populate the genotype fitness matrix and solve the system of linear equations. The following example represents a system with three divergent alleles with pathogen recognition sites according to the following layout (1 is a recognition site of a specific pathogen epitope, corresponding to a black box in figure 2.1, 0 is no recognition site of this epitope, corresponding to a white box in figure 2.1):

Allele A_1	0111111110	$w_1 = 0.8$
Allele A_2	1111111000	$w_2 = 0.7$
Allele A_3	0000000001	$w_3 = 0.1$

This system results in the following genotype fitness matrix F :

$$F = \begin{pmatrix} 0.8 & 0.9 & 0.9 \\ 0.9 & 0.7 & 0.8 \\ 0.9 & 0.8 & 0.1 \end{pmatrix} \quad (2.19)$$

These alleles have the same intrinsic merits as in the RAsOD model, and nevertheless the system is stable (solving the system of linear equations yields equilibrium frequencies for alleles A_1 , A_2 and A_3 of 65.2%, 30.4% and 4.3%, respectively). In general overdominance can maintain a k -allele polymorphism as long as the off-diagonal elements, i.e. the fitness values of the heterozygotes, are sufficiently “dominating” the diagonal elements, the fitness values of the homozygotes, and the fitness values of the heterozygotes are sufficiently close to each other, which both applies in this example. However, many other recognition site layouts with three alleles of the same intrinsic merits that do not lead to a stable polymorphism are possible in the DAA model. But is a higher allelic polymorphism in the DAA model compared to the AsOD models common or rather the exception? Chapter 3 sheds more light on this question by comparing the maintenance of alleles in various overdominance models under different scenarios.

2.4 Discussion and Conclusions

By using the general mathematical framework of a single locus multi-allele model, I classified traditional models of heterozygote advantage that do not include any spatial or temporal variation by comparing their fitness matrices. Building on a concept of balancing selection suggested by Wakeland et al. (1990) that extends traditional heterozygote advantage models by favouring differences between alleles, I developed a novel overdominance model and demonstrated how such a model fits into the general single locus multi-allele framework presented before.

Comparing all these overdominance models in the general framework clarifies the connection between the assumptions made in each of the models, and the results in terms of allele numbers, frequencies and intrinsic merits of alleles maintained at equilibrium.

For example, Lewontin et al. (1978) created random fitness matrices with no explicit overdominance model assumed. In this case, that vast majority of these fitness matrices did not represent a polymorphic equilibrium, particularly when the number of alleles increased. The only exceptions were situations that came close to symmetric overdominance, with heterozygotes having nearly equal fitness, i.e. situations that can be approximated by the genotype fitness matrix (2.12). The same fitness matrix applies to the symmetric overdominance models described by Takahata and Nei (1990) and Satta (1997). Generic asymmetric overdomi-

nance models, where heterozygotes have an equal and constant advantage over homozygotes (equation 2.13), as well as the reinforcing asymmetric overdominance model as defined by De Boer et al. (2004) with genotype fitness matrix (2.17) do not generally lead to strongly polymorphic situations at equilibrium (Lewontin et al., 1978; De Boer et al., 2004). However, other overdominance models with different genotype fitness matrices might well be both asymmetric in the fitness of the homo- and heterozygotes *and* support a larger number of alleles at equilibrium. One of these models, based on the DAA hypothesis (Wakeland et al., 1990) and with a genotype fitness matrix that depends on epitope recognition patterns, was discussed in section 2.3.4 and will be explored together with alternative overdominance models in chapter 3.

All these results and observations together help understand why different authors came to different conclusions regarding the capacity for overdominance to contribute to MHC polymorphism at equilibrium. Chapter 3 will examine the models discussed in this chapter and the amount of polymorphism they can sustain in more depth.

Chapter 3

MHC Diversity at Equilibrium under Different Overdominance Models

3.1 Introduction

Heterozygote advantage has been found to act in a wide range of vertebrate species, and particularly on MHC genes (Thursz et al., 1997; Carrington et al., 1999; Arkush et al., 2002; Stear et al., 2005). Nevertheless, the question of whether heterozygote advantage is indeed the main factor in maintaining the “extreme genetic diversity” (Apanius et al., 1997) of the MHC region has been hotly debated, and still there is no consensus (see chapter 1). Takahata and Nei (1990) and Satta (1997), using earlier theoretical work of Robertson (1962) and Karlin and Lessard (1984), argued in favour of this theory, while other authors, such as Lewontin et al. (1978) and De Boer et al. (2004) doubted an important role of heterozygote advantage in the maintenance of MHC diversity. This dispute has been further obscured by the question whether symmetric overdominance can be considered as a realistic model of heterozygote advantage, or whether the asymmetry in fitness that different alleles confer, and the asymmetry among heterozygotes, is important enough to force a move to asymmetric overdominance models.

Here I investigated the capacity of different overdominance models to maintain multiple alleles at equilibrium. All the models examined were defined in the general framework outlined in chapter 2 by their genotype fitness matrix F . I compare the results of the different overdominance models in terms of the number of persisting alleles, their frequencies and intrinsic merits, and more advanced measures of diversity such as diversity profiles ${}^qD^Z$ (Leinster and Cobbold, 2012), which will be extensively used in chapters 4, 5 and 6.

Finally, I discuss the biological realism of the models, as well as their limitations and potential model expansions.

3.2 Materials and Methods

3.2.1 Description of the allele frequency update algorithm

To explore the number of alleles n that can be maintained in a k -allele system ($n \leq k$) under different overdominance models, I used an iterative algorithm that changed the frequencies of the alleles in the gene pool from one generation to the next using equation 2.2. This is possible since an equilibrium solution of the k -allele system solves the system of linear equations given by (Karlin and Lessard, 1984)

$$(\mathbf{F}\mathbf{p}^*)_i = \bar{w}(\mathbf{p}^*) \quad i = 1, \dots, k \quad (3.1)$$

Here, $\bar{w}(\mathbf{p}^*)$ is the population fitness of the k -allele system when allele frequencies are at equilibrium. Since the left-hand side of equation 3.1 is just the marginal fitness of allele i (see equation 2.3), this implies that the marginal fitnesses are equal at equilibrium. For that reason, the method of applying allele frequency changes according to equation 2.2 in every generation ultimately leads to exactly the allele frequency distribution at equilibrium, as the marginal fitnesses become more similar from one generation to the next. The dynamics of 2.2 is ultimately governed by Fisher's fundamental theorem of natural selection (Lieberman, 1991), as $\bar{w}(\mathbf{p})$ increases from one generation to the next as long as allele frequencies are not at equilibrium.

Though the matrix formulation assumes an infinite population, the algorithm was set to remove alleles from the gene pool if their frequencies fell below a threshold value of 10^{-8} ; this was done to ensure that alleles *can* be lost. The updating of the frequencies of the remaining alleles was performed for as many generations as was necessary to ensure that the marginal fitnesses of all the alleles were equal (when rounded to 15 decimal places). This is equivalent to a state where the frequencies of the remaining alleles have reached equilibrium (\mathbf{p}^*), and hence no further allelic loss can occur.

3.2.2 Allele configurations explored

To facilitate comparison between the different overdominance models, two (arguably quite artificial) scenarios were explored for all models (see table 3.1). The first scenario started with an initial set of 100 alleles, whose intrinsic merits were uniformly distributed in the interval $[0, 1]$, with the intrinsic merits of the top and bottom allele being 0.995 and 0.005, respectively. This means that homozygotes with two top alleles had a fitness of $f_{tt} = 0.995$, which is nearly 200 times higher than homozygotes with two bottom alleles ($f_{bb} = 0.005$), and implies that individuals carrying the top allele were nearly fully protected against all

pathogens.

The second scenario started with just 50 alleles whose intrinsic merits were uniformly distributed in the interval $[0.01, 0.5]$, and the top and bottom allele had intrinsic merits of 0.01 and 0.5, respectively.

scenario	number of initial alleles	intrinsic merits of initial alleles
1	100	0.005, 0.015, 0.025, \dots , 0.995
2	50	0.01, 0.02, 0.03, \dots , 0.5

Table 3.1: Initial sets of alleles in the two scenarios explored.

These scenarios were chosen to allow for a more general interpretation of the results, as they covered situations that resulted in a low (scenario 1) or high (scenario 2) heterozygote advantage (see table 3.2). Heterozygote advantage is low for scenario 1 as only the alleles with the highest intrinsic merits ultimately persist; these are relatively close to one for scenario 1 but not for scenario 2. This results in homozygotes with a fitness close to one for scenario 1 (the fitness of heterozygotes lies between the fitness of the corresponding homozygotes and one in all overdominance models).

3.2.3 Overdominance models and diversity measures

When assessing the capacity for heterozygote advantage to maintain a stable polymorphism of k alleles, the overdominance model used plays a crucial role (see section 2.3). All overdominance models characterised in chapter 2, the symmetric overdominance (SOD) model that allows for variation in the intrinsic merit of alleles, the simple asymmetric overdominance model (SAsOD), the reinforcing asymmetric overdominance (RAsOD) model and the divergent allele advantage (DAA) model were explored. All these models except for the DAA model are fully deterministic, therefore a single run of the allele frequency update algorithm described earlier (see section 3.2.1) was sufficient to obtain the set of persisting alleles. In case of the DAA model, the recognition site patterns were drawn randomly at the start of each run of the allele frequency update algorithm, which was then applied for as many generations as was necessary to reach equilibrium.

The final set of stable alleles was used to calculate values of interest; for the DAA model, these values were averaged over all simulation runs that were performed (10000 unless specified otherwise). The overdominance models were compared in terms of two simple diversity measures, the number of alleles maintained at equilibrium and the range of intrinsic merits of these stable alleles. In special cases, more advanced diversity measures, such as diversity profiles ${}^qD^Z$ (Leinster and Cobbold, 2012), were used to compare certain models.

3.3 Results

3.3.1 Simple measures of persisting allelic diversity in the different overdominance models

Figures 3.1, 3.2 and 3.3 show the distribution of the frequencies of the persisting alleles for scenarios 1 and 2 for the SOD, SAsOD and RAsOD model, respectively. All these overdominance models had a strict relationship between the intrinsic merit of an allele and its frequency, resulting in a strict order of the allele frequencies according to their intrinsic merit. This was different for the DAA model (figure 3.4 shows the allele frequency distribution for a particular randomly drawn configuration for each scenario), where alleles with lower intrinsic merit could still exist at higher equilibrium frequencies than alleles with higher intrinsic merit, i.e. the allele frequencies were not necessarily all ordered in magnitude in the same way as the intrinsic merits of the alleles were.

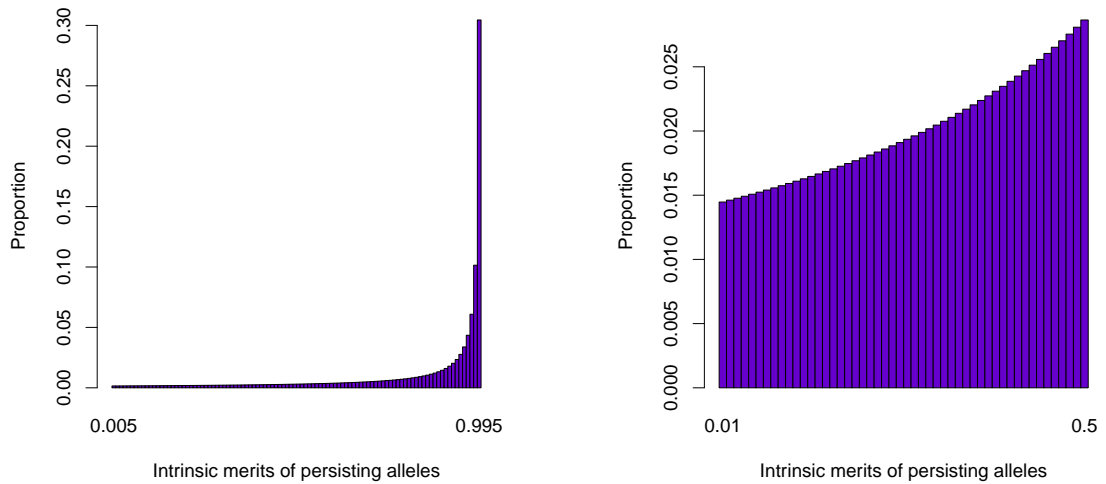


Figure 3.1: Distribution of the frequencies of alleles in the SOD model for scenario 1 (left) and scenario 2 (right).

Table 3.2 shows the number of persisting alleles (n_{equil}) and the range of intrinsic merits of the alleles that were maintained (r_{equil}), defined as the absolute difference between the allele with the highest intrinsic merit and the allele with the lowest intrinsic merit from the set of persisting alleles. It also shows the amount of heterozygote advantage $A_{het,equil}$ that was present in the final gene pool, calculated as the increase in the average fitness of heterozygotes compared to the average fitness of homozygotes. For the DAA model, the 95% confidence intervals, obtained from a Monte Carlo bootstrapping algorithm with 10000 repeats for each simulation run and the quantiles function of R (R Development Core Team, 2011), and the maxima are given for all three measures.

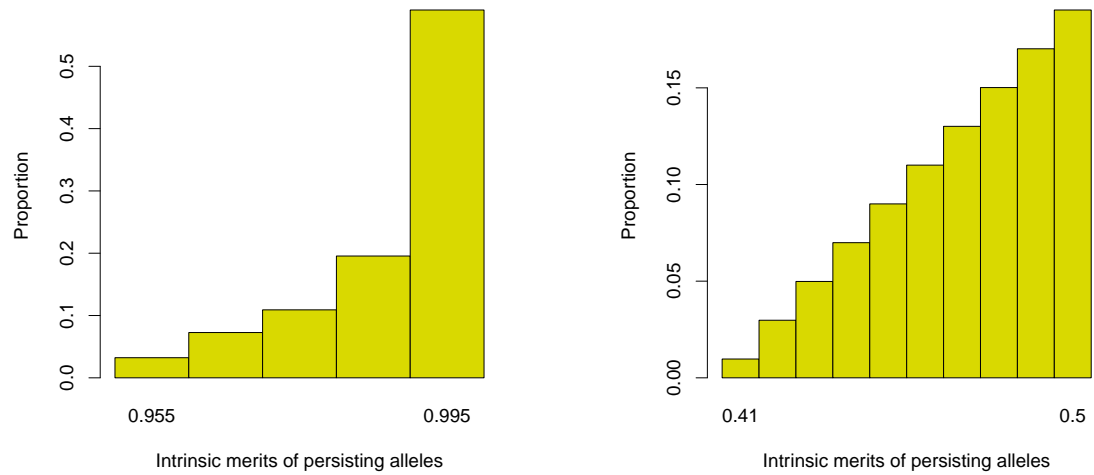


Figure 3.2: Distribution of the frequencies of alleles in the SAsOD model for scenario 1 (left) and scenario 2 (right).

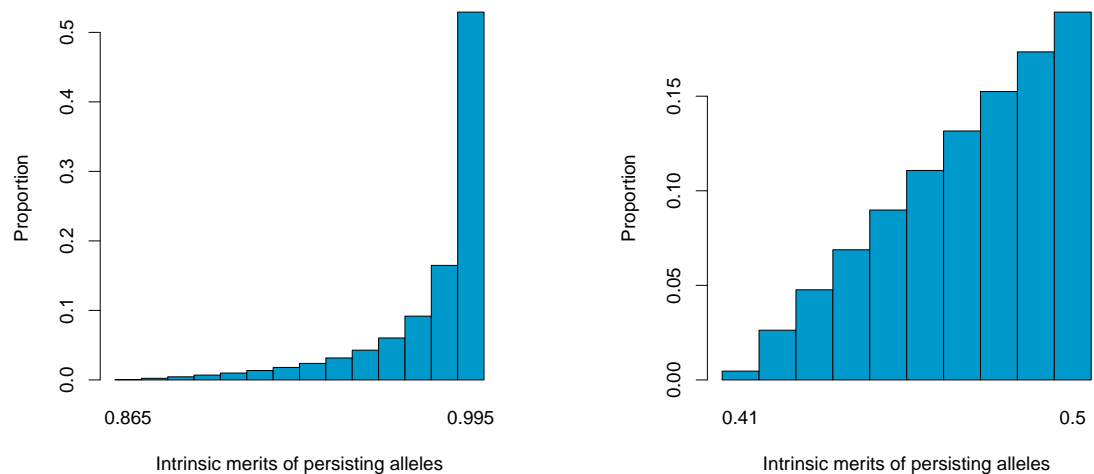


Figure 3.3: Distribution of the frequencies of alleles in the RAsOD model for scenario 1 (left) and scenario 2 (right).

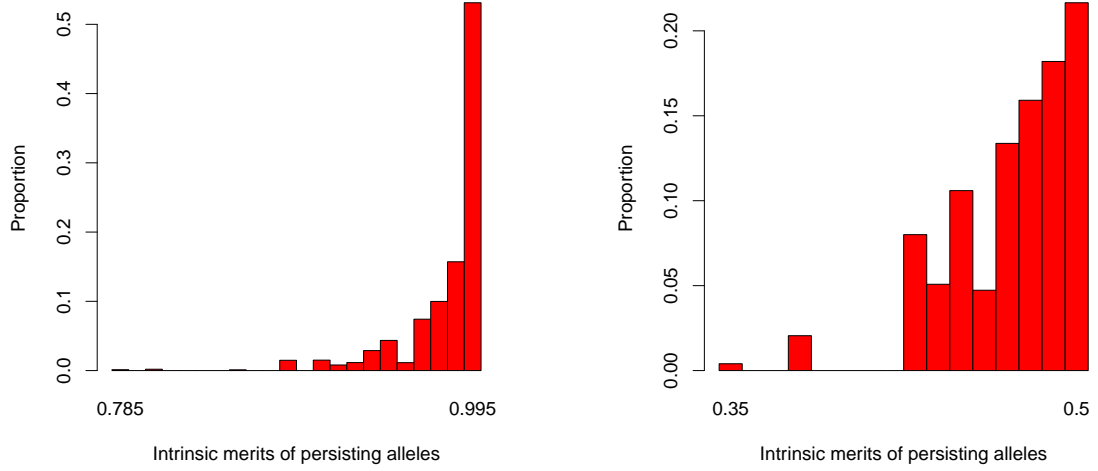


Figure 3.4: An illustrative distribution of the frequencies of alleles in the DAA model for scenario 1 (left) and scenario 2 (right). These distributions correspond to a single, random set of starting alleles.

	n_{equil}		r_{equil}		$A_{het,equil}$
	95% CI	max	95% CI	max	95% CI
scenario 1					
SOD	100	100	0.99	0.99	1.35%
SAsOD	5	5	0.04	0.04	0.69%
RAsoD	14	14	0.13	0.13	0.69%
DAA	[11, 16]	≥ 20	[0.12, 0.24]	≥ 0.25	[0.64%, 0.79%]
	n_{equil}		r_{equil}		$A_{het,equil}$
	95% CI	max	95% CI	max	95% CI
scenario 2					
SOD	50	50	0.49	0.49	220.28%
SAsOD	10	10	0.09	0.09	49.95%
RAsoD	10	10	0.09	0.09	49.95%
DAA	[7.86, 11.00]	≥ 14	[0.07, 0.12]	≥ 0.14	[48.95%, 53.86%]

Table 3.2: Number of persisting alleles (n_{equil}), range of intrinsic merits of persisting alleles (r_{equil}) and heterozygote advantage in the set of persisting alleles ($A_{het,equil}$) in different overdominance models. Only for the DAA model, the 95% confidence intervals for n_{equil} , r_{equil} and $A_{het,equil}$ are given, while these values do not vary in the other overdominance models.

The number of alleles persisting in the DAA model and the range of intrinsic merits of these alleles varied depending on how the recognition sites of the alleles A_i were set. This is demonstrated in figures 3.5 (number of alleles maintained) and 3.6 (range of intrinsic merits of persisting alleles), which show the distributions resulting from 10000 simulation runs, each of which uses a randomly drawn epitope recognition pattern for all 100 (scenario 1) or 50 (scenario 2) alleles.

For scenario 1, the number of persisting alleles in the DAA model was on average similar to the number of persisting alleles in the RAsOD model (marked in blue), but there was a fair amount of variation around this mean value, with some allele sets resulting in very low or very high numbers of persisting alleles. A similar result can be seen in scenario 2, however with a distribution that is slightly biased towards lower allele numbers.

The range of intrinsic merits of the set of persisting alleles in the DAA model also followed a distribution that is centered around the range supported in the RAsOD model for scenario 2, while the supported ranges were on average much higher in the DAA model compared to the RAsOD model for scenario 1.

The distribution of the product of the number of persisting alleles and the range of their intrinsic merits gave information about the proportion of configurations that supported large numbers of alleles while still allowing for a considerable variation in intrinsic merit of the alleles maintained (figure 3.7). For most configurations in scenario 1, this product was higher in the DAA model, while for scenario 2 the number of configurations with higher and lower products was comparable.

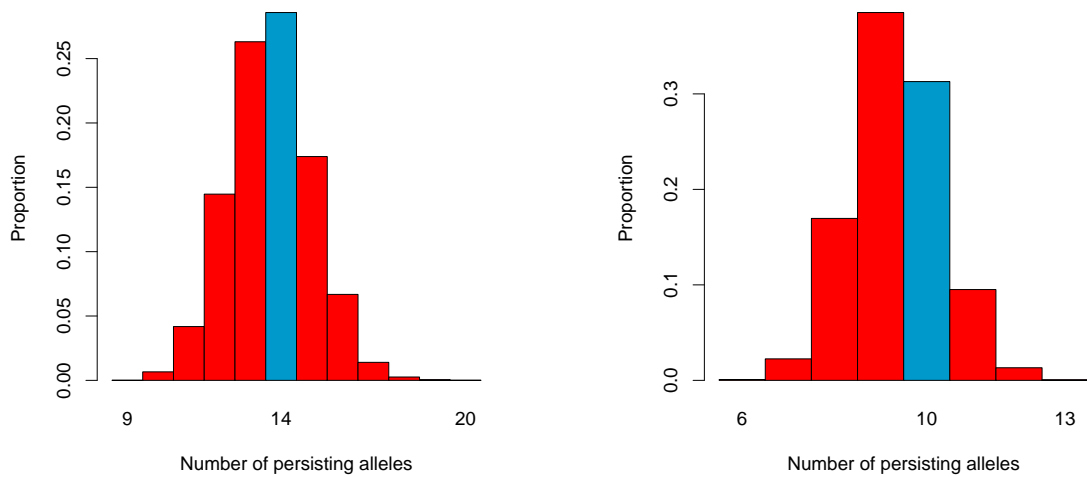


Figure 3.5: Distributions of the number of alleles maintained at equilibrium in the DAA model for scenario 1 (left) and scenario 2 (right). The blue bar marks the number of persisting alleles in the RAsOD model.

While on average the number and intrinsic merit ranges of persisting alleles were not greatly

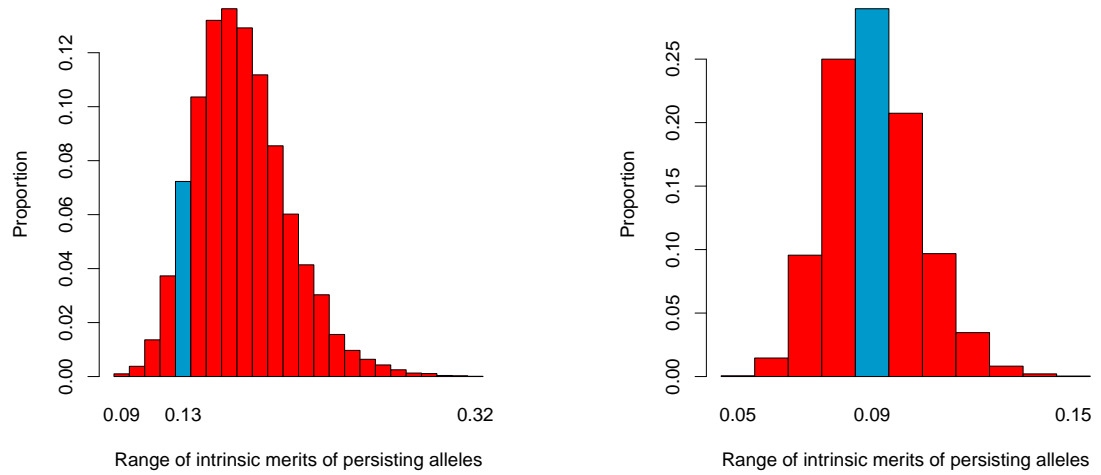


Figure 3.6: Distributions of the range of intrinsic merits of the set of alleles maintained at equilibrium in the DAA model for scenario 1 (left) and scenario 2 (right). The blue bar marks the range of intrinsic merits in the RAsOD model.

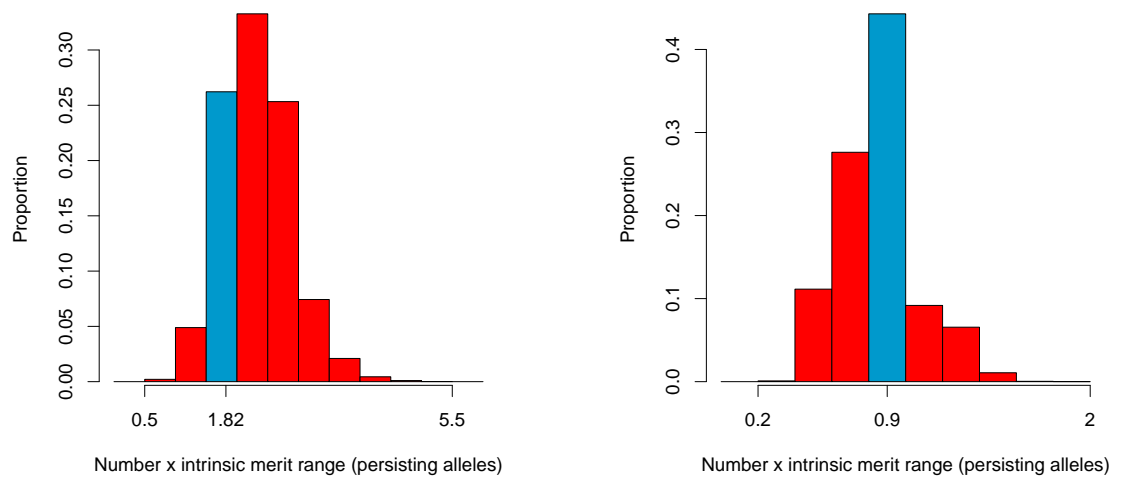


Figure 3.7: Distributions of the product of the number of alleles maintained at equilibrium and the range of intrinsic merits of these alleles in the DAA model for scenario 1 (left) and scenario 2 (right). The blue bar marks the same measure in the RAsOD model.

different between the DAA and the RAsOD models, sets of initial alleles that at the same time resulted in a higher number of persisting alleles and a larger intrinsic merit range of these persisting alleles (as compared to the RAsOD model) did exist in the DAA model. The maximum values given in table 3.2 represent the maximum found over 200000 and 1000000 simulation runs for scenarios 1 and 2, respectively. These numbers are hence lower boundaries of possible maximum values, which means that sets of alleles existed where allele numbers were at least 40% higher in the DAA than in the RAsOD model, while the range of intrinsic merits was at the same time over 40% higher for scenario 2, and over 90% higher for scenario 1, again comparing the DAA model to the RAsOD model.

3.3.2 Diversity profiles of selected overdominance models

A more sophisticated way to characterise diversity is to examine diversity profiles ${}^qD^Z$ (Leinster and Cobbold, 2012) for the different models and scenarios. Figure 3.8 compares the diversity profiles of the asymmetric overdominance models with the DAA model. The top row of figure 3.8 was generated from allele frequency data that did not take similarity between alleles into account (“naive” diversity, $Z = I$). These two plots show the diversity profiles for scenario 1 (left) and scenario 2 (right) when q varies from 0 to 10. For the figures in the bottom row a similarity measure was used, the intrinsic merit difference between pairs of alleles, since this measure can be applied to all models. Entries in the similarity matrix Z have to be normalised to the interval between 0 and 1, where 0 corresponds to two totally different alleles, and 1 corresponds to two alleles that are the same with respect to the property of interest.

When the intrinsic merit difference between pairs of alleles is used as a similarity measure, then equation 3.2 gives a possible normalisation; this normalisation is used in figure 3.8.

$$Z_{ij} = \frac{(w_{\max} - w_{\min}) - |w_i - w_j|}{w_{\max} - w_{\min}} \quad (3.2)$$

Here, w_{\max} and w_{\min} are the maximum and minimum intrinsic merits allowed (i.e. 1 and 0 in scenario 1 and 0.5 and 0 in scenario 2), and w_i and w_j are the intrinsic merits of alleles A_i and A_j . If $w_i = w_j$, then similarity between alleles A_i and A_j is 1, but if e.g. $w_i = w_{\max}$ and $w_j = w_{\min}$, then the similarity between alleles A_i and A_j is 0 (these alleles are considered entirely different).

Comparing the top row to the bottom row of figure 3.8 demonstrated that diversity increased in the DAA model relative to the asymmetric overdominance models if intrinsic merit difference was taken into account. This was a significant relative increase, as variation between the 10000 different sets of alleles was relatively low in the DAA model.

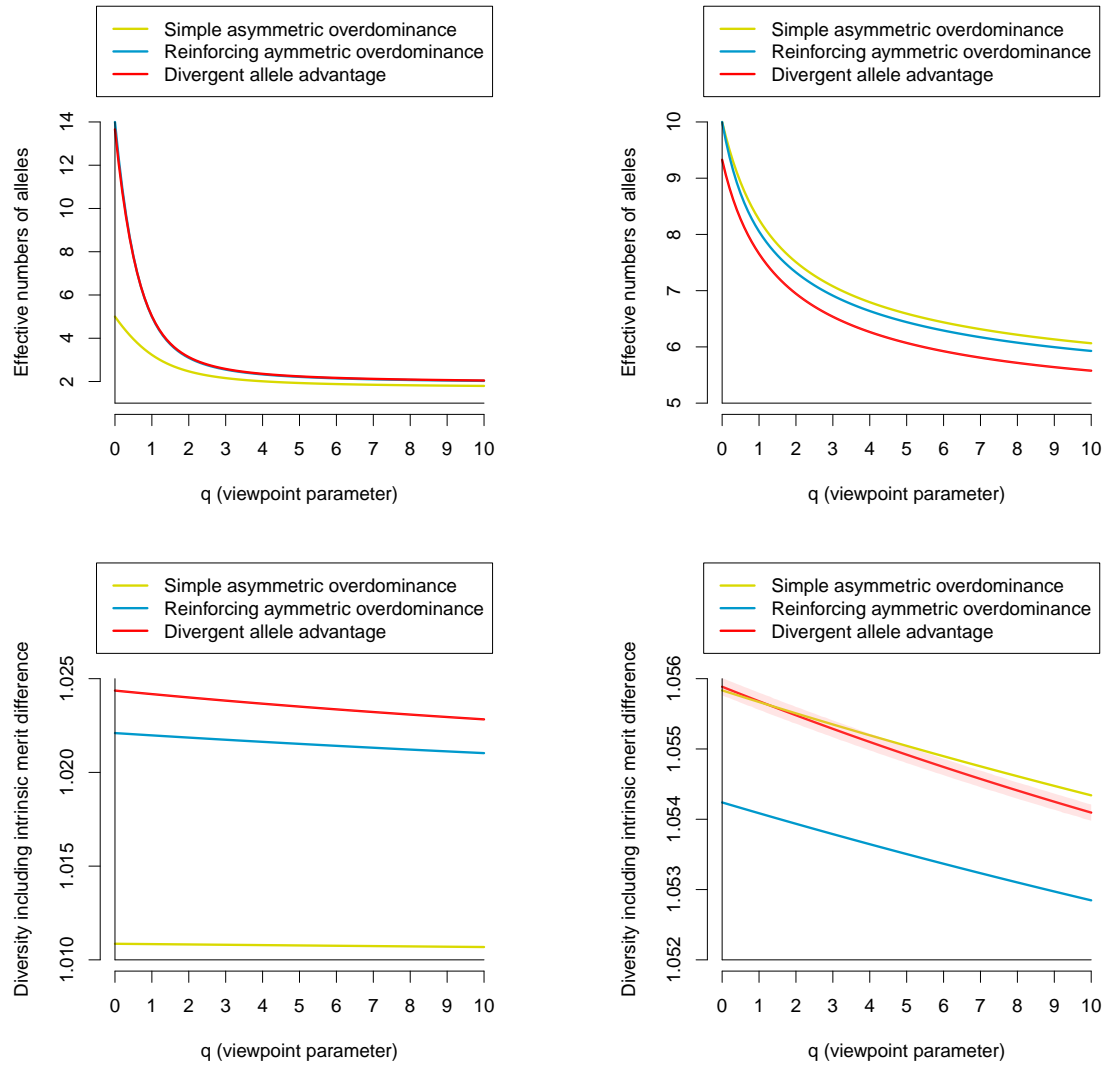


Figure 3.8: Comparison of the naive diversity profiles (top) and the diversity profiles when taking into account the similarity of alleles in terms of their intrinsic merit (bottom), for scenario 1 (left) and scenario 2 (right). The light red band shows the 95% confidence interval for the DAA model, which is too narrow to be distinguishable for all figures except the bottom right figure.

3.3.3 Diversity in the DAA model

One important question still remains: what are the differences between sets of alleles in the DAA model that result in larger numbers of persisting alleles and a higher range of their intrinsic merits, and those where the number of persisting alleles and their range is low – and which sets of alleles can be considered more diverse, in particular when using more elaborate diversity measures like diversity profiles? To answer the first part of this question, it is desirable to simplify it by examining the correlation between the number of alleles maintained at equilibrium and the range of intrinsic merits of these alleles. The setup of both scenarios, with a fixed number of starting alleles whose intrinsic merits were evenly distributed already indicates that a positive correlation between these values can be expected, and figure 3.9 confirmed this presumption. The R-squared values were 15.5% for scenario 1 and 52.3% for scenario 2, and the regression line had a distinct positive slope. For that reason only relationships involving the number of persisting alleles were subsequently investigated, while results for the range of intrinsic merits can be found in the appendix only (figures A.1 and A.2).

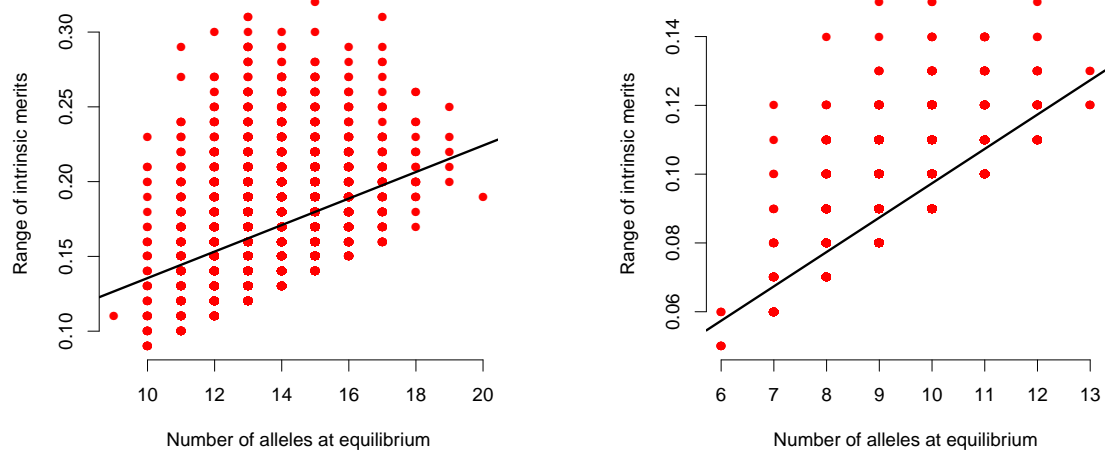


Figure 3.9: Correlation between the number of alleles maintained at equilibrium and the absolute range of intrinsic merits of these alleles for scenario 1 (left) and scenario 2 (right).

3.3.3.1 Number of persisting alleles and susceptible recognition sites of heterozygotes

The sets of alleles used when starting the simulations did not differ in intrinsic merits or frequencies of alleles they contained. The main differences between any two sets of alleles in the DAA model were the positions of “susceptible” recognition sites (white squares in figure 2.1) on the alleles. According to the definition of genotype fitness in this model (see

section 2.3.4), the proportion of susceptible recognition sites of a genotype directly translates to the fitness of this genotype. If heterozygous genotypes carry alleles that complement each other well, then these genotypes will have an advantage over the two corresponding homozygous genotypes, and the respective increased off-diagonal element in the genotype fitness matrix will stabilise the set of alleles. Also, this will result in an increase in the amount of heterozygote advantage for this allele set.

The first relationship analysed was therefore the relationship between the average proportion of susceptible recognition sites of heterozygotes and the number of alleles maintained at equilibrium. Figure 3.10 shows that there was a positive correlation between these measures, i.e. in sets where the number of persisting alleles was large, heterozygous genotypes had on average more susceptible recognition sites. The R-squared values were 51.2% and 76.3% for the unweighted allele sets (top row in figure 3.10) and scenarios 1 and 2, respectively, and 13.9% and 33.0% for the weighted allele sets (bottom row in figure 3.10) and scenarios 1 and 2, respectively.

The intuitive explanation for situations where larger numbers of alleles were maintained was that the bottom alleles (in terms of intrinsic merit) have more susceptible recognition sites, and therefore these alleles will create genotypes with higher numbers of susceptible recognition sites; however this was not the main driver of the positive relationship between the proportion of susceptible recognition sites and the number of alleles. As the frequencies of the bottom alleles were typically very low, these alleles did not have a large impact on the average proportion of susceptible recognition sites once genotypes were weighted by their frequencies (bottom row in figure 3.10), which means that such an explanation did not generally apply. An alternative explanation emerges when taking only the m most frequent alleles of each situation into account. This approach avoided any bias from different numbers of alleles in the sets of persisting alleles, and comparisons of average proportions of susceptible recognition sites and their relationship to the number of alleles maintained were more meaningful.

Figure 3.11 shows the same relationship as figure 3.10, except that the average proportion of susceptible recognition sites was calculated from the $m = 8$ (in case of scenario 1) and $m = 4$ (in case of scenario 2) most frequent alleles only. Again there was a pronounced positive correlation.

Interestingly, the same positive correlation existed for all m ($3 \leq m \leq n_{min}$) if m is the number of top alleles and $n_{min} = 10$ for scenario 1 and $n_{min} = 8$ for scenario 2 (these values for n_{min} were chosen in order to include all situations where m was less than or equal to the number of persisting alleles for the majority of repetitions). This is shown in figure 3.12. The relationship between the two measures had a distinct positive slope that was similar for all m ; this was true for both scenario 1 and scenario 2, and it applied whether

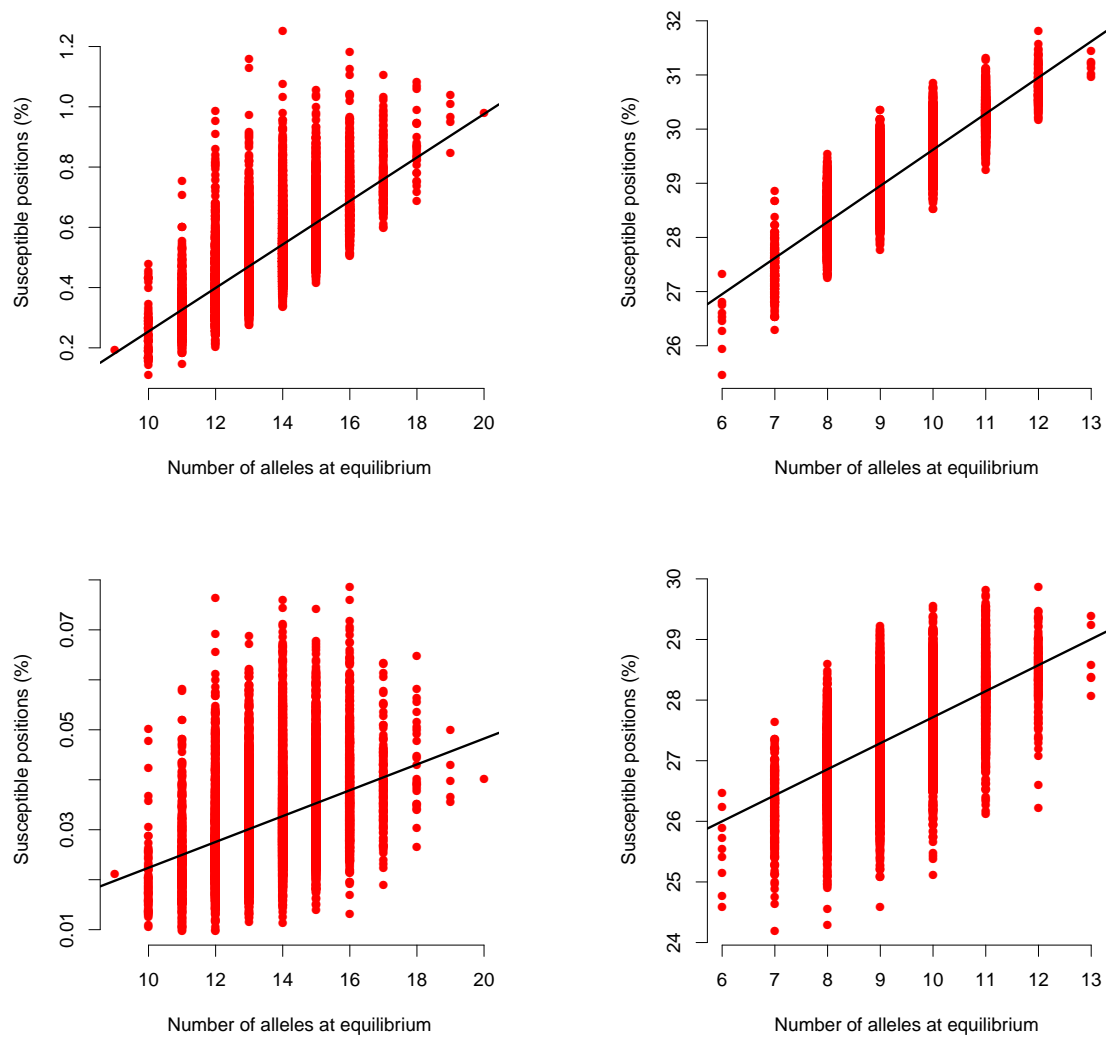


Figure 3.10: Average proportion of susceptible recognition sites of heterozygotes vs. number of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). In the top row all persisting alleles are weighted equally, whereas in the bottom row the proportions of the heterozygous genotypes are used as weights. A positive correlation can be seen for all situations.

genotypes were weighted according to their equilibrium frequencies (bottom row of figure 3.12) or equally weighted (top row of figure 3.12).

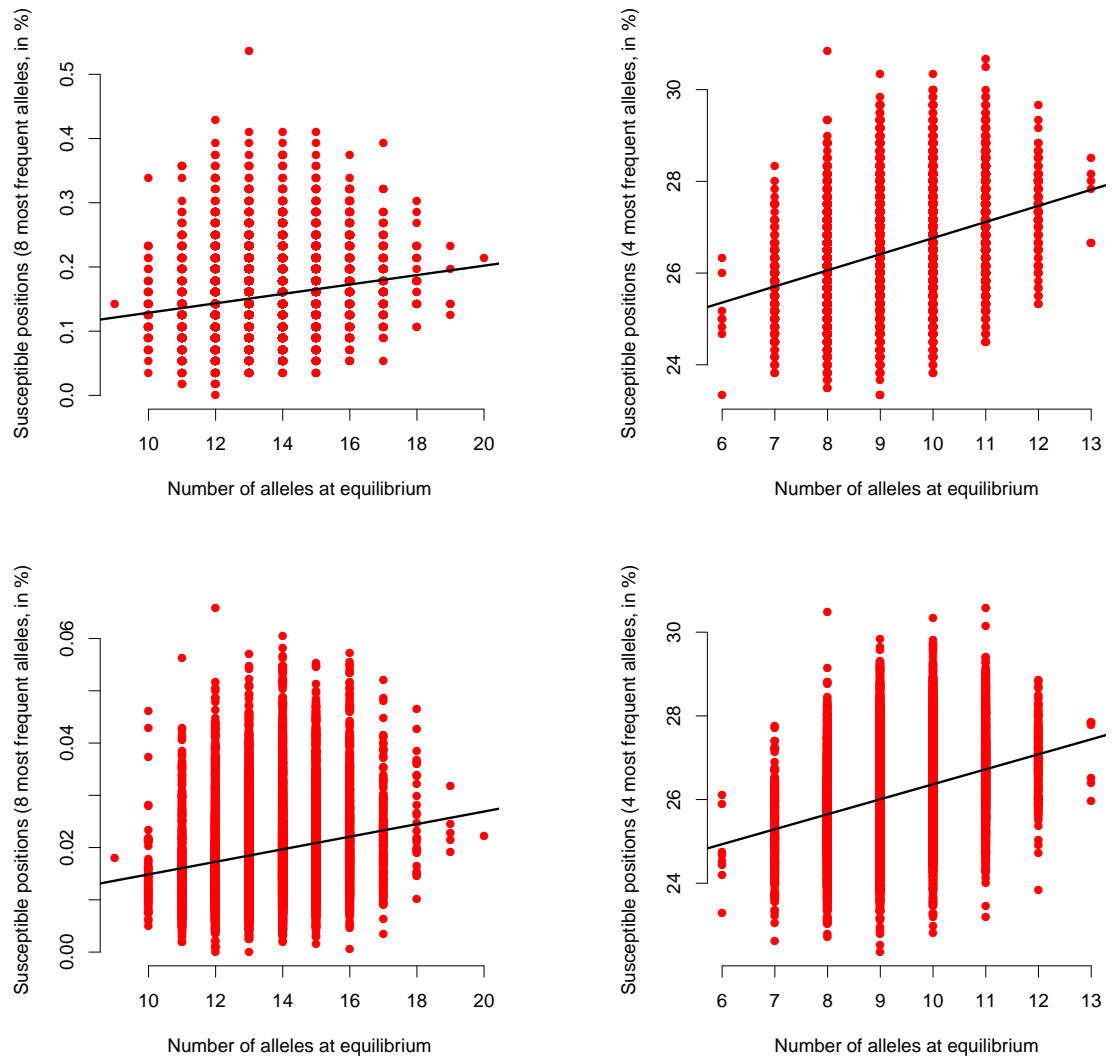


Figure 3.11: Average proportion of susceptible recognition sites of heterozygotes vs. number of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). Only the eight (scenario 1) and four (scenario 2) most frequent alleles were taken into account. For the two top figures, all maintained alleles had the same weight, whereas the two bottom figures use the equilibrium frequencies of the heterozygous genotypes as weights. A positive correlation can be seen for all situations.

3.3.3.2 Number of persisting alleles and population fitness

Another interesting aspect is that there was a relationship between the number of persisting alleles and the population fitness. These measures showed a pronounced negative correlation (figure 3.13), which was a somewhat surprising result, given that intuitively more alleles are expected to increase heterozygosity and thus population fitness in an overdominance system.

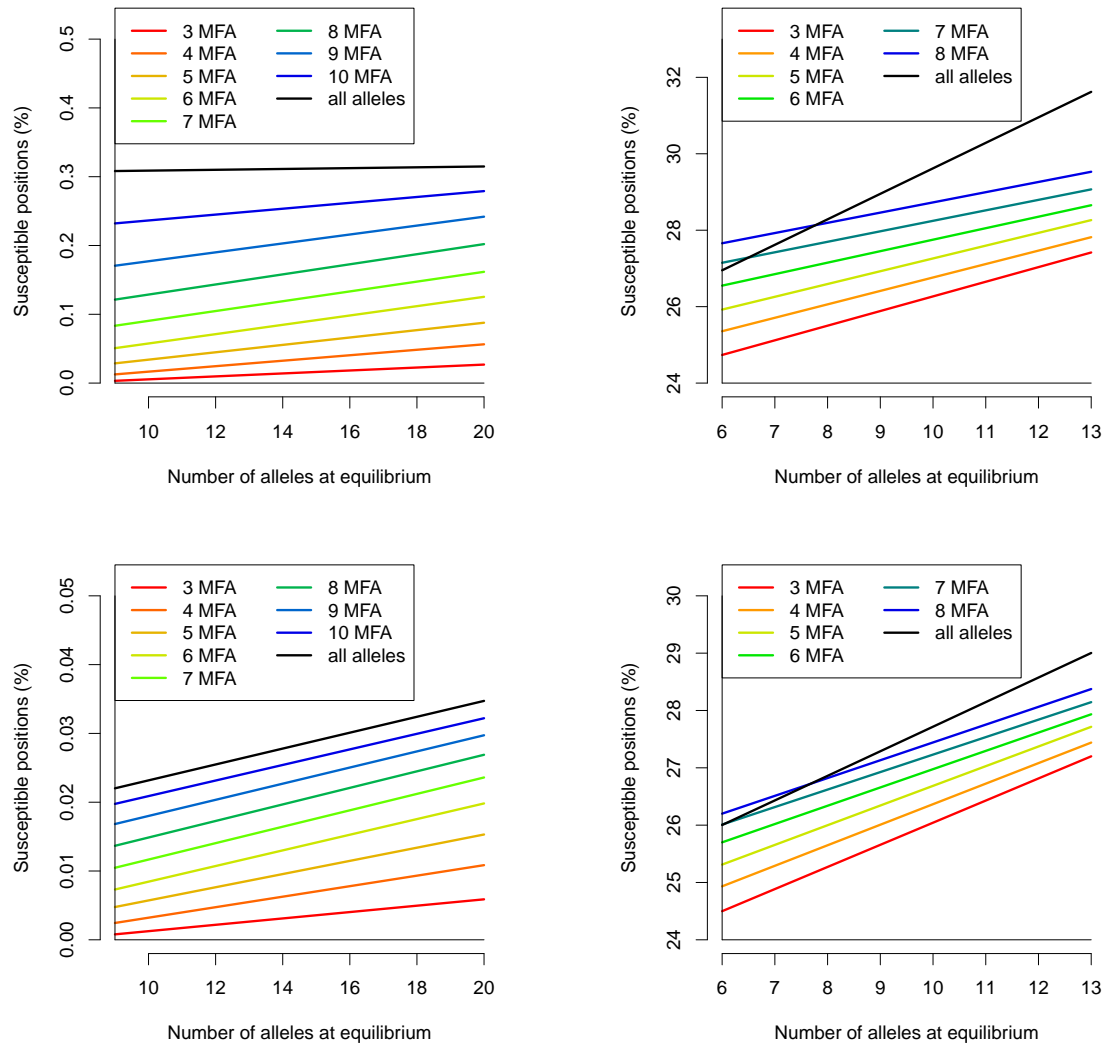


Figure 3.12: Average proportion of susceptible recognition sites of heterozygotes vs. number of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). Only the j most frequent alleles (MFA) were taken into account, with j varying from 3 to 10 for scenario 1 and 3 to 6 for scenario 2. The black regression line corresponds to all alleles taken into account. For the two top figures, all maintained alleles had the same weight, whereas the two bottom figures use the equilibrium frequencies of the heterozygous genotypes as weights. A positive correlation can be seen for all situations.

In this case, however, this result was just an implication of the positive correlation between the average proportion of susceptible recognition sites of heterozygotes and the number of alleles maintained at equilibrium. If heterozygotes have more susceptible recognition sites, then population fitness decreases – this was confirmed by figure 3.14, which demonstrates a relatively strong negative correlation between these measures, with an R-squared value of 48.4% for scenario 1 and 85.6% for scenario 2.

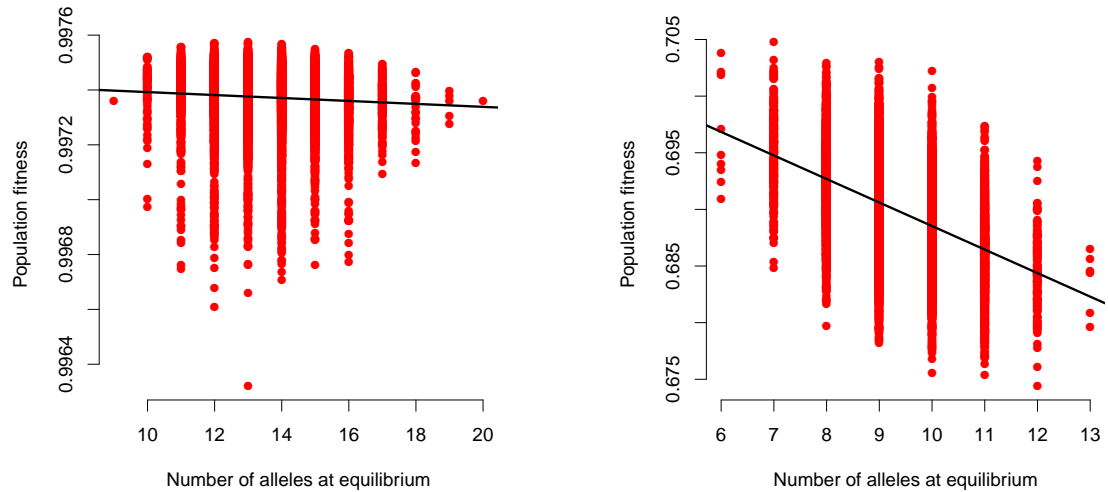


Figure 3.13: Population fitness vs. number of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). A negative correlation can be seen for both scenarios.

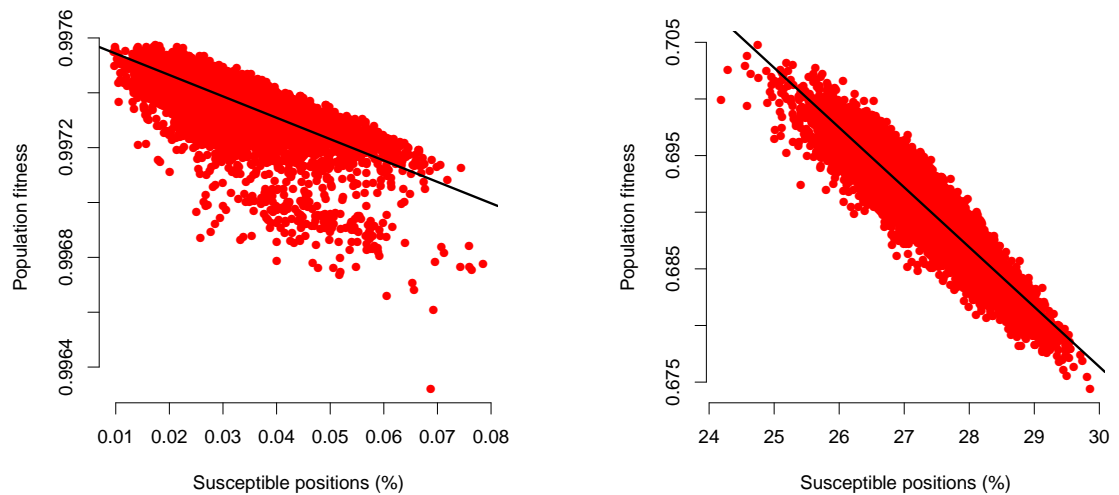


Figure 3.14: Population fitness vs. average proportion of susceptible recognition sites of heterozygotes for scenario 1 (left) and scenario 2 (right). Both figures use the proportions of the heterozygous genotypes as weights. A negative correlation can be seen for both scenarios.

3.3.3.3 Diversity profiles for different sets of alleles in the DAA model

While it was demonstrated that lower numbers of alleles persisting at equilibrium were correlated with a better covering of the epitope recognition space by the most frequent alleles, the question still remains whether these sets of alleles are in fact less diverse than sets that maintain a higher number of alleles, but where the most frequent alleles cover the epitope recognition space less well. Clearly, an answer to this question depends on how diversity is defined.

I approached this question by using diversity profiles (Leinster and Cobbold, 2012) as a measure of diversity, both for naive diversity that does not take similarity measures into account, and for similarity-sensitive diversity. I used two different similarity measures: in the first case I defined dissimilarity as the difference between intrinsic merits of alleles (as in section 3.3.2); for the second case I defined dissimilarity as the proportion of differences in the epitope recognition sites:

$$Z_{ij} = \frac{l_t - l_{ij}}{l_t} \quad (3.3)$$

Here, l_t denotes the total number of epitope recognition sites of an allele (200 in scenario 1 and 100 in scenario 2), while l_{ij} is the number of differences in the epitope recognition sites between alleles A_i and A_j . This similarity measure is used since it is strongly connected to the notion of coverage of the epitope recognition space by genotypes (see below), and therefore, as shown in figure 3.14, to the population fitness.

In terms of both their naive and similarity-sensitive diversity profiles it is necessary to discuss scenarios 1 and 2 separately, as the implications are different. I subsequently examine diversity profiles for numbers of persisting alleles ranging from 10 to 18 in case of scenario 1 and 7 to 12 for scenario 2, as only these cases were supported by a large enough number of allele sets (see figure 3.5).

For scenario 1, the naive diversity profiles for numbers of persisting alleles ranging from 10 to 18 intersected, and therefore figure 3.15 splits this set of diversity profiles into three viewpoint parameter intervals for ease of viewing. For low values of the viewpoint parameter q , i.e. if more emphasis is put on the absolute number of alleles compared to the evenness of allele frequencies, diversity was highest for sets where the number of persisting alleles was high, and steadily decreased to sets where the number of persisting alleles was low (top left figure). This changed however when q increased, and more emphasis was put on the evenness of allele frequencies. The profile curves crossed between $q = 1$ and $q = 2$ (top right figure). For the highest values of q , diversity was greatest for sets where the number of persisting alleles was low, and steadily decreased to sets where the number of persisting alleles was high (bottom left figure), in contrast to the situation for small q . This implies that in scenario 1, an answer to the question “which sets of alleles exhibit on average a greater

naive diversity, sets with a larger or a smaller number of persisting alleles”, depends on the viewpoint, reflected in the viewpoint parameter q . For scenario 2, the answer is simpler; here, sets with a higher number of persisting alleles also exhibited a larger naive diversity for the whole range of the viewpoint parameter q (figure 3.15, bottom right).

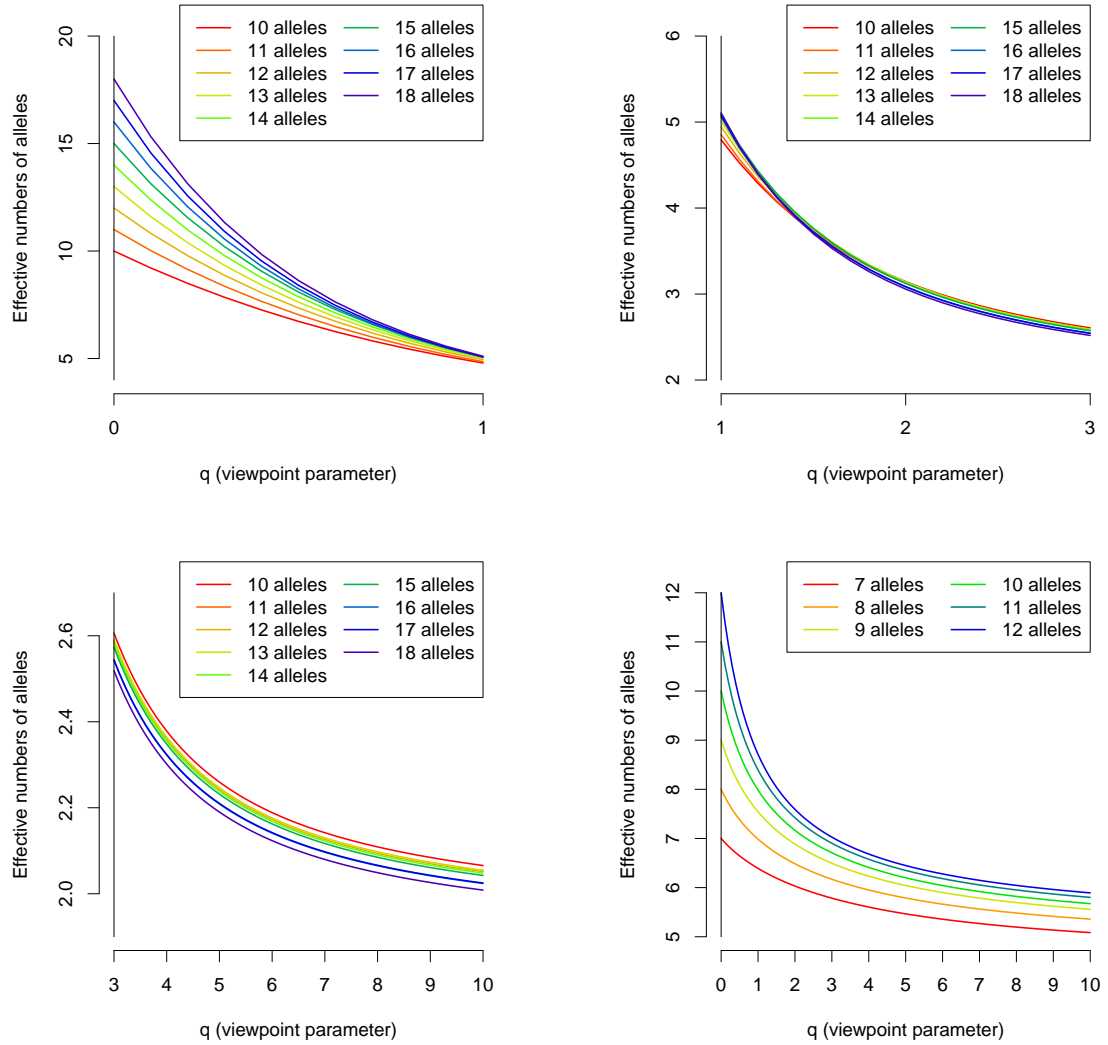


Figure 3.15: Diversity profiles for naive diversity for scenario 1 (top left, top right and bottom left) and scenario 2 (bottom right). The naive diversity profile for scenario 1 is split into three viewpoint parameter intervals for ease of viewing. The coloured lines correspond to sets with different numbers of persisting alleles.

When using the intrinsic merit as a measure for allele similarity, just as in the overdominance model comparison of section 3.3.2, then sets with a higher number of persisting alleles were more diverse for all values of q in both scenarios (figure 3.16). This is in line with figure 3.9 and the strong correlation between the number of alleles maintained and the range of their intrinsic merits. In scenario 1, this effect was strong enough to reverse the order of curves in the naive diversity profile for larger values of q .

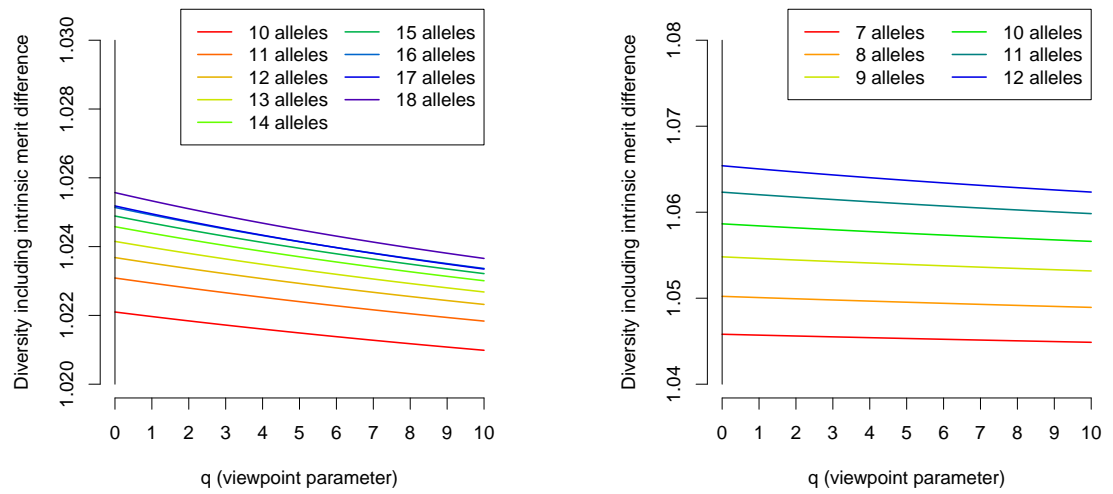


Figure 3.16: Diversity profiles for similarity-sensitive diversity for scenario 1 (left) and scenario 2 (right). The coloured lines correspond to sets with different numbers of persisting alleles. This similarity-sensitive diversity profile uses differences between intrinsic merit of alleles as a measure for similarity between alleles.

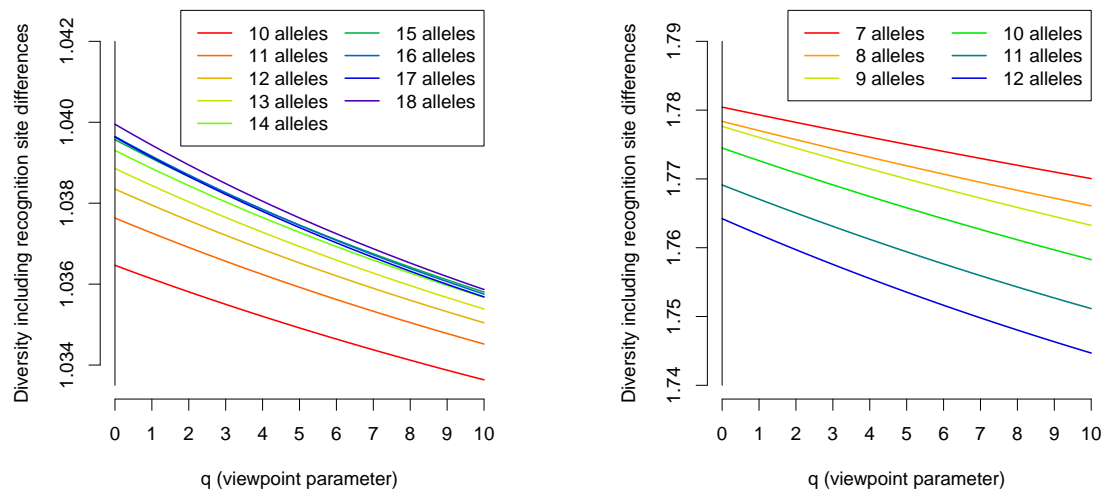


Figure 3.17: Diversity profiles for similarity-sensitive diversity for scenario 1 (left) and scenario 2 (right). The coloured lines correspond to sets with different numbers of persisting alleles. This similarity-sensitive diversity profile uses differences between epitope recognition sites of alleles as a measure for similarity between alleles.

Another possible measure of allelic diversity specifically for the DAA model discussed here is the proportion of differences in the epitope recognition sites between alleles of heterozygotes (figure 3.17 shows the respective diversity profiles). For this measure the two scenarios behaved differently. In scenario 1, there was a positive correlation between the differences in the epitope recognition sites of heterozygotes and the susceptible positions of heterozygotes, both weighted by genotype frequency (R-squared value of 21.0%, figure 3.18 (left)), whereas in scenario 2, the same correlation was negative (R-squared value of 91.1%, figure 3.18 (right)).

When examining the diversity profiles themselves (figure 3.17), then the diversity when using the differences in the epitope recognition sites between alleles as a similarity measure decreased from sets with large numbers of alleles maintained to sets with small numbers of persisting alleles for scenario 1 and all q (therefore reversing the order of the naive diversities for larger values of q). For scenario 2, however, the opposite is true, and diversity increased from large to small numbers of persisting alleles, again for all q (reversing the order of the naive diversities for all q).

These results are now particularly interesting in the context of population fitness. In scenario 1, the sets with the highest numbers of persisting alleles were the most diverse sets when including epitope recognition site differences (figure 3.17, left). Genotypes consisting of alleles of these sets had, however, a higher average proportion of susceptible sites (figure 3.18, left), which in turn resulted in a lower population fitness (figure 3.14, left), in line with earlier results (3.13, left).

Scenario 2 behaved differently, as sets with the lowest numbers of persisting alleles were the most diverse sets when including epitope recognition site differences in this scenario. However, as the heterozygotes consisting of alleles of sets with the lowest numbers of persisting alleles (i.e. sets of alleles where the top alleles complement each other particularly well) had a lower average proportion of susceptible sites (figure 3.18, right), population fitness nevertheless decreased for sets with large numbers of persisting alleles, in line with figure 3.13 (right). This is because population fitness was negatively correlated with the average proportion of susceptible sites of heterozygotes (figure 3.14, right). Therefore, in scenario 2, sets with low numbers of persisting alleles were the most diverse ones when including recognition site differences, and these sets had a low average proportion of susceptible sites of heterozygotes, and thus a high population fitness.

3.4 Discussion

One of the primary aims of any theory attempting to explain the extraordinary polymorphism at the MHC is to explain large numbers of alleles. Classical theories can be broadly

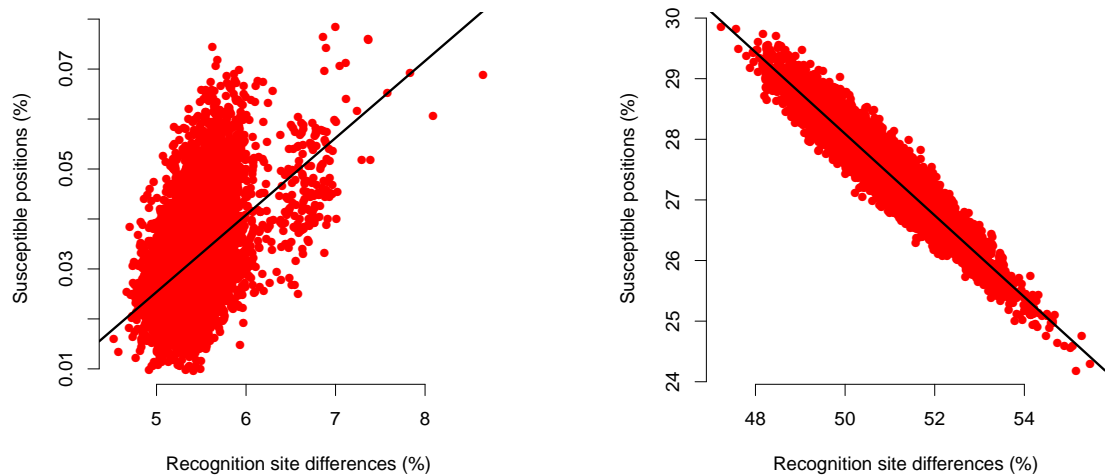


Figure 3.18: Average proportion of susceptible recognition sites of heterozygotes vs. average proportion of recognition site differences of heterozygotes, for scenario 1 (left) and scenario 2 (right).

grouped into two categories, symmetric overdominance models, where all the heterozygotes have the same fitness value, and asymmetric overdominance models, where the fitnesses of homozygotes as well as heterozygotes may vary, but heterozygotes are on average fitter than homozygotes.

While symmetric overdominance models are often preferred when a theory capable of explaining large numbers of alleles is required (Takahata and Nei, 1990; Stoffels and Spencer, 2008), their biological realism is sometimes disputed (De Boer et al., 2004), as there is some evidence of considerable fitness differences between heterozygotes (Bronson et al., 2013). However, the more asymmetric overdominance models are, the lower the degree of MHC diversity they are able to explain (reviewed in Hedrick (1999)). These restrictions have led authors to reject heterozygote advantage as a primary driver of MHC diversity altogether (Lewontin et al., 1978; De Boer et al., 2004).

The divergent allele advantage (DAA) model, built on the concept of balancing selection suggested by Wakeland et al. (1990) and discussed in chapter 2, could explain the MHC polymorphism better, as it is more closely linked to the biology of the MHC, and intrinsically allows for asymmetry in the fitness of heterozygotes. This made it desirable to include the DAA model in the equilibrium analysis performed in this chapter.

In the symmetric overdominance (SOD) model explored here, all alleles of the starting set were maintained in both scenarios, therefore the range of intrinsic merits of alleles was always the entire range of the set of initial alleles. This is confirmed by earlier theoretical and simulation results, which demonstrated that symmetric overdominance is a mechanism which is capable of maintaining large numbers of alleles (Takahata and Nei, 1990; De Boer

et al., 2004). However, such a model struggles to explain allelic loss other than from stochastic events, and would therefore predict the persistence of alleles that confer a very low fitness on homozygotes at equilibrium.

In both the simple (SAsOD) and the reinforcing asymmetric overdominance (RAsOD) models, only the alleles with the highest intrinsic merits persist at equilibrium. This is because alleles with low intrinsic merit create homozygotes, but also heterozygotes with low fitness, and therefore the marginal fitness values of these alleles stay below the population fitness, causing a decrease of the frequencies of these alleles in every generation until they fall below the threshold value and are removed from the population.

The DAA model behaves similarly to the asymmetric overdominance models, with the main difference that allele frequencies at equilibrium are not necessarily arranged according to size, and “holes” where alleles with higher intrinsic merit are absent, while alleles with lower intrinsic merit are maintained (see figure 3.4), can exist. This is important, as it demonstrates that to a degree alleles can compensate low intrinsic merit by possessing epitope recognition patterns that complement the most frequent alleles well.

Comparing traditional asymmetric overdominance models to the DAA model, it was shown that the DAA model can lead to a higher or lower amount of polymorphism for the same set of initial alleles.

If the polymorphism is measured in terms of the number of alleles maintained, then the average value for the DAA model is similar to the number of persisting alleles in the RAsOD model for both scenario 1 and 2. When comparing the range of the intrinsic merits of the persisting alleles, then the DAA model is on average more polymorphic than the RAsOD model for scenario 1, but exhibits a similar level of diversity for scenario 2. The same is true if the product of the number of alleles maintained and the range of the intrinsic merits is considered. The latter measure is important, as the maintenance of more alleles requires a narrower window for the intrinsic merits of the persisting alleles (De Boer et al., 2004). This result alleviates concerns about the capacity for overdominance in general to maintain both larger numbers of alleles and variation in intrinsic merit of these alleles (De Boer et al., 2004), as a substantial proportion of simulations in both scenarios showed a greater allele-range product in the DAA model compared to the RAsOD model.

When comparing diversity for a range of different viewpoints by using similarity-sensitive diversity profiles together with intrinsic merit of alleles as the similarity measure, then the DAA model is more diverse than both asymmetric overdominance models for scenario 1, and more diverse than the RAsOD model (and similar in diversity to the SAsOD model) for scenario 2.

The amount of polymorphism when using simple diversity measures, such as the number of alleles maintained or the range of the intrinsic merits of these alleles, ultimately depends

on the particular epitope recognition patterns of the alleles in the initial set. The results demonstrated that if the most frequent alleles (which are in general also the alleles with the highest intrinsic merits) complement each other well, i.e. if the genotypes emerging from this pool of complementary alleles have on average only few susceptible recognition sites, then there is little room for additional alleles, and the number of alleles maintained at equilibrium is low while the population fitness is high. If, however, the most frequent alleles do not complement each other particularly well, then more additional alleles can coexist with these top alleles at equilibrium, particularly if these alleles can contribute recognition sites that the top alleles lack (i.e. if their recognition site profile complements the more frequent alleles well).

This is true if the genotypes are weighted according to their expected frequency, but, importantly, also if the genotypes are all equally weighted. Therefore knowledge of the recognition patterns of the most frequent alleles already provides information about the potential for a smaller or larger number of persisting alleles, and a smaller or larger range of intrinsic merits of these alleles, without any knowledge about their equilibrium frequencies.

Moreover, knowledge of the average proportion of differences between alleles of heterozygotes and the correlation of this value with the average proportion of susceptible positions of heterozygotes can already provide a wealth of information that can be used to infer diversities of configurations – this was demonstrated when assessing whether sets of alleles with lower or higher numbers of persisting alleles are more diverse – and the population fitness relative to other configurations.

Finally, this chapter demonstrated that a detailed analysis is necessary already for very simple setups, as scenarios 1 and 2 sometimes yielded unexpected or contrary results.

3.5 Conclusions

For the examined setup, SOD models can maintain by far the greatest amount of diversity, as all alleles will be maintained at equilibrium, independent of their intrinsic merits. The DAA model can support a moderately greater amount of diversity than traditional asymmetric overdominance models. This is true not only for well-established single diversity measures like number and range of persisting alleles, but also for the diversity of order q for all q , $0 \leq q \leq \infty$, when using the family of diversity measures ${}^qD^Z$ (Leinster and Cobbold, 2012) and a similarity matrix Z based on the intrinsic merit difference between alleles.

Another important result is the significance of the epitope recognition patterns of the most frequent alleles in the DAA model. The diversity of the set of persisting alleles, and subsequently the population fitness, is mainly determined by the degree to which the recognition patterns of these frequent alleles can complement each other.

3.5.1 Model limitations

All the aforementioned models consider only single locus multi-allele systems. The MHC however is a multigene family, and some areas of the MHC are preferentially inherited as haplotypes. Although reducing this complex region of the genome to a single locus might seem like an oversimplification, the main properties can be expected to still apply when intrinsic merits are assigned to haplotypes instead of genes. As linkage disequilibrium is strong between genes of the MHC class II region (Sommer, 2005), such a redefinition of a locus is not unreasonable. However, expanding the models to multiple loci, and examining how this affects the MHC polymorphism would be an interesting line of future research.

Mutation is another phenomenon that is not explicitly considered in these models. This is because only internal stability of a k -allele system was discussed. When looking at external stability of such a system, and the invasion prospects of a newly mutated allele, then mutation needs to be considered as well.

Genetic drift is clearly only captured if the population size is finite, as is, for example, the case in an agent-based model. Such a model can assess the strength of genetic drift compared to the strength of selection. Another aspect that is important in real-world populations is the question whether the population is even at equilibrium. Since most environments are in dynamic change, a stable equilibrium may be unlikely to ever be reached. A slow approach to equilibrium, together with permanent environmental change, could result in higher numbers of alleles present at any point in time, even to the extent that under certain assumptions, these fluctuations are able to maintain a polymorphism of the MHC without any overdominance present (Hedrick, 2002).

3.5.2 Perspectives

This chapter demonstrated that the DAA model has the potential to maintain higher numbers of alleles at equilibrium and at the same time a wider range of intrinsic merits compared to the RAsOD model. If selection favours divergence between alleles, resulting in a reduced overlap in epitope recognition sites and therefore in an improved coverage of the epitope recognition space, then the DAA model is capable of explaining an even larger allelic diversity. This in particular means that the assertion that heterozygote advantage cannot explain observed diversities of MHC genes, as argued in De Boer et al. (2004), is only true for the particular asymmetric overdominance model used in this work, but does not have general validity.

Of all the overdominance models compared here, the DAA model is the only model that has a mechanistic biological foundation for variation in genotype fitnesses. In addition to that, the equilibrium analysis has shown that this model can explain a larger allelic diversity compared

to traditional asymmetric overdominance models under certain conditions, in particular if selective forces drive sequence divergence. Furthermore, there is empirical evidence that this form of balancing selection might act in particular host-pathogen systems (Wakeland et al., 1990; Lenz et al., 2009; O'Connor et al., 2010; Lenz et al., 2013). All this suggests that the DAA model, based on the divergent allele advantage hypothesis, has the potential to be a powerful mechanism maintaining diversity at the MHC.

However, to provide a more definitive answer to whether this mechanism might even be the fundamental force maintaining the polymorphism at the MHC, finite population sizes, mutations and genetic drift need to be taken into account. I present a more sophisticated model based on the DAA hypothesis in chapter 4, which is scrutinised in a setting that takes all the evolutionary forces into account.

Chapter 4

The Maintenance of MHC Diversity in a Finite Population Subject to Evolutionary Forces

4.1 Preface

This chapter was compiled from a draft manuscript of an article that is intended for submission to a scientific journal. Therefore, the introduction as well as the methods section review some of the material covered in previous chapters.

4.2 Introduction

One of the main aims of modern genomics is to identify the genes that influence susceptibility to disease. For Mendelian traits, the presence or absence of disease is determined entirely or predominantly by the alleles at a single locus. Most infectious, parasitic and autoimmune diseases are more complex; susceptibility is determined by a number of environmental factors and multiple genes. Hundreds of studies in thousands of individuals have shown that the distribution of gene effects is L-shaped. A small number of genes have relatively large effects while the majority of genes that are associated with disease have only a small influence on the development of disease (Strachan and Read, 2010).

Surprisingly little effort has gone into analysis of how specific alleles influence disease, in particular how alleles interact with other alleles at the same locus. Yet this information is crucial for determining which individuals are susceptible. For example, a genotype with a single copy of a recessive susceptibility allele has a different prognosis than a genotype with two copies of the same allele. With heterozygote advantage, the influence of an allele also

depends upon the other allele at the locus. A lack of understanding of these interactions is particularly severe for the major histocompatibility complex which is the most important region of the vertebrate genome in determining susceptibility to disease, but currently only single allele effects on relatively rare autoimmune diseases are used to predict resistance to disease (Browning and McMichael, 1996). The purpose of this chapter is to determine how MHC alleles interact to influence susceptibility to infectious and parasitic diseases. This is one of the most pressing problems in disease genetics. In addition, the approach may serve as a model for post genomic analysis of other loci.

In most species of higher vertebrates, the MHC is the most diverse region of the genome (Klein, 1986; Browning and McMichael, 1996). In humans, over 2000 alleles have been identified at the class I HLA-A, B and C loci, and over 1000 at the class II HLA-DRB1 locus (EMBL-EBI, 2015). Fewer individuals have been studied in other species but even so tens to hundreds of alleles have been identified (Sommer, 2005). In contrast, several species such as moose and bison exhibit very low levels of MHC diversity (Mikko et al., 1999).

The MHC codes for antigen-presenting proteins that determine the specificity and intensity of the immune response and therefore plays a critical role in resistance against disease (Browning and McMichael, 1996), but the precise mechanisms maintaining the high levels of polymorphism remain unknown. The failure to understand how the MHC contributes to variation in resistance to disease prevents us from optimally managing the health of humans, livestock and wildlife populations.

The immunological function of the antigen-presenting proteins suggests that heterozygotes will recognise a wider variety of parasite molecules (“parasite” includes both macroparasites and microbial pathogens). Therefore heterozygote advantage has been suggested as the primary driving force that maintains MHC diversity (Doherty and Zinkernagel, 1975). Studies in humans (Thursz et al., 1997; Carrington et al., 1999), sheep (Stear et al., 2005), fish (Arkush et al., 2002) and wildlife (Oliver et al., 2009) have demonstrated that heterozygotes are more resistant to specific diseases than homozygotes. However, genetic models demonstrate that traditional forms of heterozygote advantage (symmetric and asymmetric overdominance models) cannot maintain large numbers of alleles unless all alleles confer very similar fitness (Lewontin et al., 1978; De Boer et al., 2004). The predictions of these traditional models are contradicted by the existence of a large number of studies that associate particular alleles with disease (Radwan et al., 2010), leading to rejection of the predominant role of heterozygote advantage and the widespread belief that resistance to one disease confers susceptibility to other diseases (Doherty and Zinkernagel, 1975; Apanius et al., 1997). This has prevented us from understanding and exploiting MHC variation in selective breeding of managed populations and hindered us from identifying humans with increased susceptibility to disease.

Here, I reconcile the contrasting genetic and immunological viewpoints by examining a mechanism of intra-locus allelic interaction called divergent allele advantage (DAA) (Wakefield et al., 1990). This mechanism is a special form of heterozygote advantage. Under the DAA hypothesis, heterozygotes with divergent alleles recognise a wider set of parasite molecules than heterozygotes with similar alleles or homozygotes, implying more effective immune responses against each parasite and also against a wider range of parasites. Under both the DAA model and traditional asymmetric overdominance models, the marginal fitness of an allele depends on its intrinsic merit, defined as the fitness of a homozygote carrying that allele, and its interaction with other alleles. The critical advance over traditional models of heterozygote advantage is that in the DAA model the fitness of a heterozygote is also influenced by the number of amino acids that differ between the two alleles.

DAA is being explored experimentally (Richman et al., 2001; Eizaguirre et al., 2012; Lenz et al., 2013), but distinguishing between competing hypotheses of parasite-mediated selection from observational data remains challenging (Spurgin and Richardson, 2010). An alternative approach is to simulate the evolution of MHC diversity under DAA. This approach has not been previously possible because when a new species is created by evolution, it will usually contain an unknown number of MHC polymorphisms. However, both sheep and cattle share an inversion in their MHC, which means that all extant MHC sequences in these two subfamilies stem from a common ancestral sequence that predates their divergence 20–40 million years ago (Lewin et al., 1999; Decker et al., 2009). I used this unique feature of the bovidae to underpin a model of the evolution of MHC diversity and rigorously assess whether this special form of heterozygote advantage explains observed diversity at the MHC.

4.3 Materials and Methods

4.3.1 Model overview

I ran stochastic simulations over a period of 10 million non-overlapping generations. This corresponds to approximately 20–40 million years, which reflects the split of the bovidae from the cervidae (Jermann et al., 1995; Decker et al., 2009). An inversion in the MHC implies that all extant MHC loci in the bovidae descend from a single ancestral chromosome (Lewin et al., 1999).

The evolution of a single gene was modelled, with an allele at that locus represented by the amino acid sequence it encodes. Only alleles with different amino acid sequences were viewed as distinct alleles. The default number of variable positions on the protein chain encoded by an allele was $n = 82$, which is approximately the length of exon 2 of the class II antigen-presenting loci. This exon contains the antigen binding site and most polymorphisms

occur in this exon. I used the simplifying assumption that each of the 20 amino acids is equally likely at any of the n positions of an allele (Brown et al., 1993; Lenz, 2011).

The models presented in this chapter differ from the DAA-inspired model presented in chapters 2 and 3 in terms of how an allele is represented. The DAA-inspired model of chapters 2 and 3 regarded an allele as a distinct pattern of pathogen peptide recognition sites, whereas a more traditional notion of an allele being represented by the encoded amino acid sequence is used in this chapter.

4.3.2 Changes in allele frequencies per generation

Frequency changes of the alleles present in generation t were modelled using the allele frequency update algorithm presented in section 3.2.1 together with a procedure that ensured that allele frequencies remained discrete multiples of the frequency of a single allele representative. As an allele representative I consider a single copy of an allele in the gene pool of a population. If for example just two alleles A and B are present in a population of 10000 individuals, then the gene pool of this population might consist of 15000 representatives of allele A and 5000 representatives of allele B (the number of representatives of both alleles A and B sums to 20000, the total number of allele copies in the gene pool).

The allele frequencies in generation $t + 1$, \mathbf{p}_{t+1} , (where bold indicates a vector) depended on the frequencies in generation t , \mathbf{p}_t , (Lewontin et al., 1978) via

$$(\mathbf{p}_{t+1})_i = (\mathbf{p}_t)_i \cdot \frac{(\mathbf{w}_m(\mathbf{p}_t))_i}{\bar{w}(\mathbf{p}_t)} \quad (4.1)$$

where \mathbf{w}_m is the vector of marginal fitnesses of the alleles, which depends on their intrinsic merit and contribution to heterozygotes, given by

$$(\mathbf{w}_m(\mathbf{p}_t))_i = \sum_{j=1}^k f_{ij}(\mathbf{p}_t)_j \quad (4.2)$$

where f_{ij} is the fitness of a genotype with alleles i and j . The population fitness \bar{w} , the weighted average of the genotype's fitnesses, is given by

$$\bar{w}(\mathbf{p}_t) = \sum_{j=1}^k (\mathbf{p}_t)_j \cdot (\mathbf{w}_m(\mathbf{p}_t))_j \quad (4.3)$$

Given a population size m , allele frequencies had to be multiples of $q_{\min} = \frac{1}{2m}$. To ensure whole-numbered multiples of q_{\min} , I drew new allele frequencies each generation from a multinomial distribution; this allowed for stochastic loss of alleles through genetic drift, as opposed to earlier models (Spencer and Marks, 1992).

In this simulation step selection and drift are handled at the same time by drawing from a multinomial distribution with the average next-generation frequencies as weights; it precedes the mutation step (see section 4.3.4) of the simulation in every generation.

4.3.3 Selection under different overdominance models

Differences in fitness among heterozygotes and homozygotes and their respective frequencies determined the marginal fitnesses of the alleles. The fitness of a homozygous individual, f_{ii} , was given by the intrinsic merit of the allele it carries, w_i , while the fitness of a heterozygous individual depended on the overdominance model used. In case of the traditional symmetric and asymmetric overdominance models (equations 4.4, 4.5 and 4.6), the amino acid sequences that are encoded by the two alleles a heterozygote individual carries have no impact on the fitness of this individual, while under the DAA model (equation 4.7) the sequences do impact fitness, as divergent sequences will more strongly increase the individual's fitness. For all the models, however, amino acid sequences are generated (and may mutate) in exactly the same way, and intrinsic merits are assigned to newly created alleles represented by these sequences according to the same procedure (see also section 4.3.4).

In a symmetric overdominance model, all heterozygotes have the same fitness by definition:

$$f_{ij} = d \quad (\forall i \neq j) \quad (4.4)$$

Here, d is a constant and $d \geq f_{ii}$ for all i ensures overdominance. I denote the special case of an overdominance model where the fitness of all homozygotes is also equal, i.e. $f_{ii} = \tilde{d}$ with $\tilde{d} < d$, as “fully symmetric” or a fully symmetric overdominance (FSOD) model, while the more general case that allows for differences between the intrinsic merit of alleles (w_i) is just denoted as symmetric overdominance (SOD), even if it is asymmetric in the fitnesses of the homozygotes.

In the simple asymmetric overdominance (SAsOD) model, the fitness of a heterozygote (f_{ij}) is defined as

$$f_{ij} = \frac{1}{2}(w_i + w_j) + c \quad (4.5)$$

where w_i and w_j are the intrinsic merits of alleles i and j , and c is a constant. This model assumes additive effects of the two alleles, plus an additional, constant fitness contribution for heterozygotes.

In addition, an alternative model of asymmetric overdominance was examined, subsequently called reinforcing asymmetric overdominance (RAsOD), which was based on De Boer et al. (2004). In its simplest form, this model gives the fitness of a heterozygote (f_{ij}) as

$$f_{ij} = w_i + (1 - w_i) \cdot w_j \quad (4.6)$$

As the degree of heterozygote advantage is an emergent property of this model, it is not a free parameter. Consequently, this alternative model could only be compared against the other overdominance models for settings 2 and 3 (see section 4.3.5 for a definition of the settings); for these settings both asymmetric overdominance models (the SAsOD and the RAsOD model) gave essentially identical results. I therefore restrict discussion to comparisons between the SAsOD model and the DAA model in many cases.

In the DAA model, f_{ij} is defined as

$$f_{ij} = \max(w_i, w_j) + k \cdot \delta \quad (4.7)$$

where w_i, w_j are the intrinsic merits of alleles i and j ; k is the number of amino acid differences between the two alleles, and δ is the fitness benefit per amino acid difference. Here, the fitness of a heterozygote depended not only upon the intrinsic merits of its two alleles, but also on the number of amino acid differences between their sequences (Wakeland et al., 1990), each of which adds a fixed contribution to the fitness of the heterozygote. This sequence divergence based overdominance contribution relates this model even more directly to the DAA hypothesis than the DAA-inspired model presented in chapters 2 and 3. Although both these models are called the “DAA model” as they convey the same concept, they are nevertheless in fact two different models.

The overdominance contribution based on sequence divergence accounts for a sequence-dependent heterozygote advantage in a similar way as a model presented by Satta (1997), but the model applied here is distinguished by additionally accounting for differences between the intrinsic merit of alleles. Furthermore, using the same principle of an allele-based contribution (allele intrinsic merit) and a contribution of intra-locus interactions for obtaining the fitnesses of heterozygotes for all overdominance models is a major improvement over models that draw the fitness of heterozygotes randomly from some distribution (Marks and Spencer, 1991) or lack the component for intra-locus interactions (Spencer and Marks, 1992), particularly when aiming for comparing different overdominance models.

4.3.4 Mutations

Each generation, non-synonymous point mutations may occur at any position in the amino acid sequence with equal probability. Since any site could affect protein structure or function, I allowed for mutations at all sites. If an allele existed previously, i.e. if it was a back-mutation, then its intrinsic merit was set equal to that of the earlier allele, otherwise, its intrinsic merit was drawn from a Gaussian distribution with a mean \bar{w}_{new} given by

$$\bar{w}_{new} = \frac{w_{orig} + w_{init}}{2} \quad (4.8)$$

where w_{orig} is the intrinsic merit of the originator allele (the allele that gave rise to the new mutation) and w_{init} the intrinsic merit of the initial allele (the allele present at the start of the simulation). The standard deviation of the Gaussian distribution was σ , the variation in intrinsic merit of alleles.

New alleles were generated solely by point mutation with a probability of μ per codon per generation. The number of point mutations per allele in a generation was drawn from a binomial distribution, and the positions of the amino acid sequence where point mutations occur were drawn without replacement. All amino acids other than the current one were assumed to be equally likely replacements. Although gene conversion can increase diversity in the MHC (Ohta, 1995; Parham and Ohta, 1996), there is a lack of information on the frequency and extent of gene conversion; I therefore chose not to include gene conversion in this model. This is a conservative assumption, as it reduced the number of new alleles that were generated.

The mutation step succeeds the selection and drift step (section 4.3.2) and is ultimately followed by a recalculation of the marginal fitnesses of the alleles and the population fitness at the end of each generation.

4.3.5 Parameter space explored by evolutionary simulations

The key parameters were the mutation rate, μ ; the size of the well-mixed population, m ; the number of amino acids per allele, n ; the variation in intrinsic merit of alleles, σ ; and the relative heterozygote advantage per amino acid difference, δ . The ranges for these parameters are summarised in table 4.1. By varying the last two parameters (σ and δ) six different “settings”, i.e. (σ , δ)-pairs, were defined. One of these settings, subsequently called setting 1, was used as the standard setting that generated the main results. Simulations with parameter values assigned to the other five settings were then performed to assess the robustness of these results.

The mutation rate, μ , was given a value of $1.22 \cdot 10^{-8}$ per codon, in line with previous estimates (Takahata and Nei, 1990; Kumar and Subramanian, 2002; Borghans et al., 2004; Stoffels and Spencer, 2008); to assess the robustness of my results I varied this parameter from $6.1 \cdot 10^{-9}$ to $2.44 \cdot 10^{-8}$. The number of variable positions (n) on the protein chain encoded by an allele was given a default value of 82; I varied this parameter from 41 to 164.

The population size, m , varied between 10 thousand and 10 million. Estimates for population sizes of bovid species can reach hundreds of millions, given historical accounts of herd sizes of bison alone (Berger and Cunningham, 1995). Estimates of historical effective population sizes of bovid species are in the range of 10000 – 100000 for *Bos taurus* cattle (Gautier et al., 2007; de Roos et al., 2008; Goddard and Hayes, 2009; MacEachern et al., 2009) and up to

Parameter	Range of values (standard value)	Parameter description
μ	$6.1 \cdot 10^{-9} - 2.44 \cdot 10^{-8}$ ($1.22 \cdot 10^{-8}$)	nonsynonymous point mutation probability per codon per generation
m	10000 – 10 million	size of the well-mixed population
n	41 – 164 (82)	number of amino acids per allele
σ	0.006 - 0.08 (0.04)	variation in intrinsic merit of alleles
δ	0.06% – 1.17% (0.39%)	relative heterozygote advantage per amino acid difference

Table 4.1: Parameter ranges used for the simulations. Default values for the parameters (setting 1) in brackets.

250000 for the Beringia bison (Shapiro et al., 2004; Drummond et al., 2005). Therefore, the interval between 10 thousand and 10 million individuals provided a reasonable range for population size in my models.

The remaining parameters, σ and δ , were chosen to give a range of relative heterozygote advantage, Δ_r , at the end of the simulations. The range encompassed observed data and allowed exploration of settings with extreme ratios of heterozygote advantage to variation in allele intrinsic merit (table 4.2). Specifically, for a given pair of values of σ and δ I calculated the heterozygote advantage, Δ_r , (measured relative to the average homozygote fitness) at the end of the simulation, as follows:

$$\Delta_{rel} = \frac{\bar{f}_{het} - \bar{f}_{hom}}{\bar{f}_{hom}} \quad (4.9)$$

Estimates for Δ_r were based on published estimates of the selection coefficient, s , acting on MHC genes using:

$$\Delta_{rel} = \frac{s}{1 - s} \quad (4.10)$$

Estimates for s in humans vary with the MHC locus and range from less than 0.05 (Hill et al., 1991; Satta et al., 1994) to 0.3 - 0.46 (Hedrick, 1994; Black and Hedrick, 1997). Settings 1, 2 and 3 correspond to selection coefficients, s , from 0.19 to 0.42 which translate to values for Δ_r from 23% to 71% (see table 4.2). The settings with low heterozygote advantage, i.e.

settings 4, 5 and 6, correspond to a selection coefficient of 0.04, which translates to a value of 4.2% for Δ_r (see table 4.2). The settings were different in the choice of σ , as theory (Lewontin et al., 1978) and exploratory simulations showed the ratio of Δ_r to σ influenced the simulation outcomes. Setting 1 has an intermediate Δ_r to σ ratio and occupies a central position in both the $\sigma - \delta$ and $\sigma - \Delta_r$ parameter space (figure 4.1). Settings 2 and 3 have higher Δ_r to σ ratios. Settings 4, 5 and 6 have low heterozygote advantage. Setting 4 has a ratio of Δ_r to σ similar to settings 1 and 2. Settings 5 and 6 are extreme settings with low heterozygote advantage relative to variation in the intrinsic merit of alleles, with setting 6 being the most extreme and already similar to pure dominance.

Setting 1	Setting 2	Setting 3	Setting 4	Setting 5	Setting 6
(σ, δ)	(σ, δ)	(σ, δ)	(σ, δ)	(σ, δ)	(σ, δ)
Δ_r	Δ_r	Δ_r	Δ_r	Δ_r	Δ_r
(0.04, 0.39%)	(0.08, 1.17%)	(0.02, 1.02%)	(0.006, 0.06%)	(0.02, 0.07%)	(0.08, 0.10%)
23%	59%	71%	4.2%	4.2%	4.2%

Table 4.2: Model parameters controlling selection. The (σ, δ) -pairs and Δ_r for the 6 settings examined. The parameter σ is the variation in intrinsic merit of alleles, relative to the intrinsic merit of the initial allele; δ is the increase in fitness of the heterozygote individual for each difference in the MHC allele sequences, relative to the intrinsic merit of the initial allele; and Δ_r is the extent of heterozygote advantage in the population (measured relative to the average homozygote fitness).

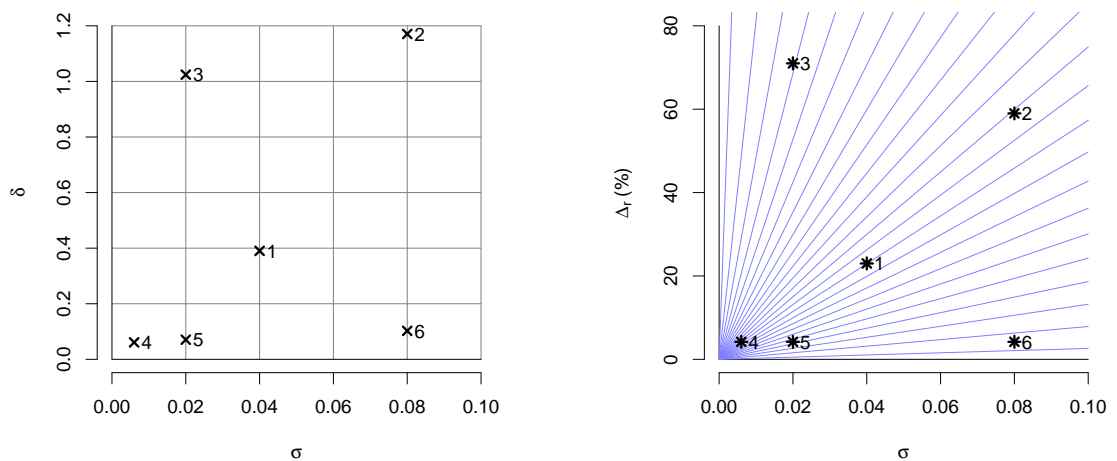


Figure 4.1: The positions of the (σ, δ) - and (σ, Δ_r) -pairs of the 6 settings examined (see table 4.2) are marked. Settings were chosen in order to cover situations with different ratios of heterozygote advantage to the variation in the intrinsic merit of alleles.

I ran simulations using the DAA model and the SAsOD model for all settings. The RAsOD

model, for reasons given in section 4.3.3, was only simulated for settings 2 and 3, and the SOD model was only simulated for settings 2 and 5, representing situations with a high and low ratio of Δ_r to σ .

The sensitivity to the choice of setting was explored for several key outcomes: the number of distinct alleles in the population (figures 4.7, 4.8 and B.5); the persisting variation in allele intrinsic merit at the end of the simulation (figures 4.11, B.11, B.12, B.13, B.14, B.15, B.16, B.17 and B.18 show the distribution of allele intrinsic merits, whereas figures 4.12 and B.19 show the weighted standard deviation in allele intrinsic merit); the expected heterozygosities (figures 4.15, B.30 and B.51a); the gene pool shares of the five most frequent alleles (figures 4.16 and B.31); the diversity profiles $^qD^I$ (figures 4.23, B.50 and B.51b); and the distribution of the emergence times of alleles (figures 4.24, B.52, B.53, B.54, B.55, B.56, B.57, B.58 and B.59). In each case, the DAA model was compared to other overdominance models, most prominently the SAsOD model.

4.3.5.1 Heterozygote advantage

I adjusted the free parameter determining heterozygote advantage in each model (δ in the DAA model, c in the SAsOD model and d in the SOD model) so that the emergent Δ_r matched either the heterozygote advantage of the RAsOD model (settings 2 and 3) or a target value (settings 1, 4, 5 and 6). In each case this was done for a population size of 1 million individuals. As the SOD model generated large numbers of alleles and was therefore computationally intensive, simulations running up to the full time span of 10 million generations were only possible for population sizes up to 250000, while simulations for a population size of 1 million were stopped after 1 million generations. This shorter time span resulted in an overestimation of the amount of heterozygote advantage, since the average fitness of homozygotes is increasing over time due to fitter alleles replacing less fit alleles. For that reason, the data point corresponding to a population size of 1 million was not taken into account, but the amount of heterozygote advantage at a population size of 1 million was rather predicted by linear regression (see figures 4.2 and B.1). It is therefore rather this extrapolated value that matches the heterozygote advantage of the other models. Figures 4.2 and B.1 demonstrate that the amount of heterozygote advantage for a population size of 1 million at the end of 10 million generations is similar for all overdominance models simulated in the respective setting.

Particularly for the DAA model, the heterozygote advantage may – again depending on the population size – take a long while to converge on the final level. Since this model directly links heterozygote advantage to the evolution of sequence divergence, the latter needs to develop first. In this model, selective forces drive alleles to be different from each other. This is shown in figures 4.3, B.2 and B.3. This feature has major implications for population fitness

recovery after genetic bottlenecks of populations (section 5.3.3 discusses this observation and the implications it has in more detail).

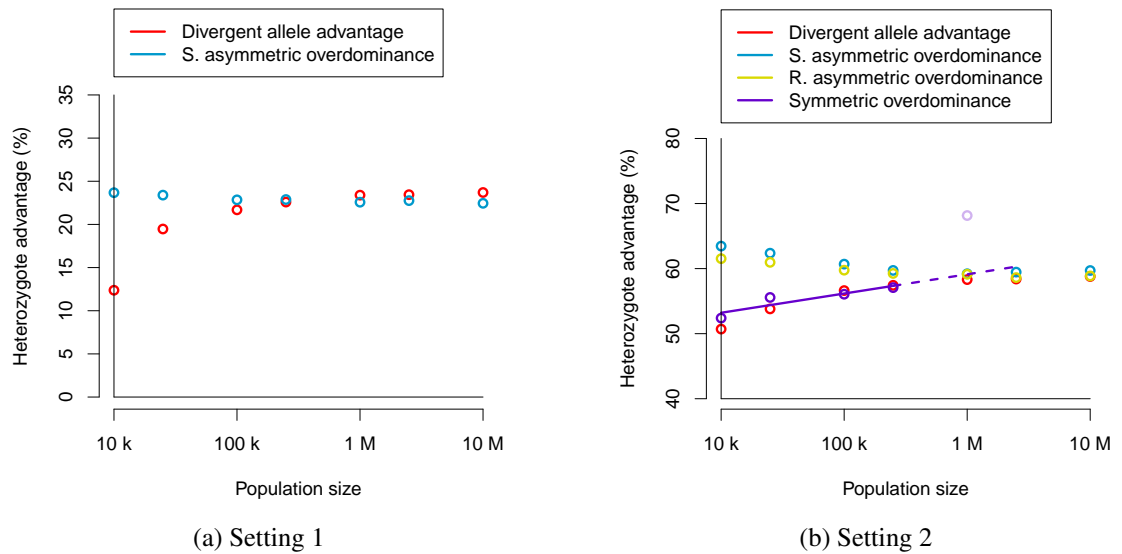


Figure 4.2: Heterozygote advantage for different overdominance models, settings 1 and 2 and various population sizes.

4.3.6 Stochastic simulations

The population size remained constant throughout the simulation. Because of inherent stochastic variation, I ran a number of repeats for each simulation. The number of repeats chosen depended on population size - the larger the population, the smaller the stochastic effects and the smaller the number of repeats. The number of repeats for both the DAA and the asymmetric overdominance models was 50, 40, 20, 10, 10, 5, 3 for population sizes of 10000, 25000, 100000, 250000, 1 million, 2.5 million and 10 million, respectively. The only exception was setting 3, where the high computational demands required the number of repetitions to be reduced to 5 for 1 million and 3 for 2.5 million. Similarly, computational demands made it necessary to reduce the number of repeats in the SOD model to 20, 10, 5, 5 and 5 for population sizes of 10000, 25000, 100000, 250000 and 1 million, respectively, and for a population size of 1 million, the number of generations was reduced from 10 million to 1 million. Larger population sizes were not simulated in the SOD model.

4.3.7 Comparison with MHC diversity in natural populations

The observed data comes from ten studies on frequencies of the DRB allele in bovidae populations, which are summarised in table 4.3.

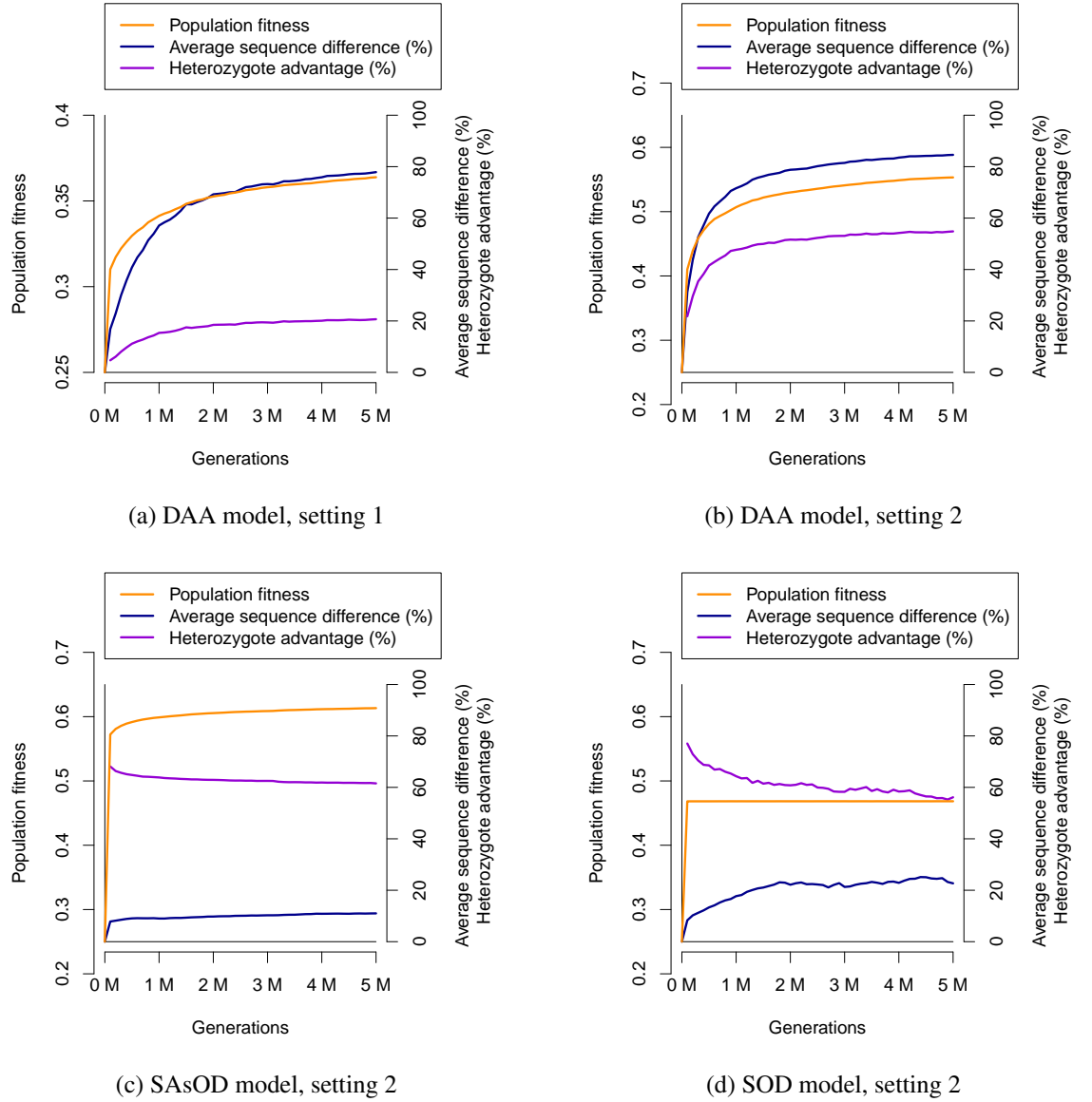


Figure 4.3: Average sequence difference, heterozygote advantage and population fitness for different overdominance models, settings 1 and 2 and a population size of 100000. The approach to population fitness equilibrium is much slower in the DAA model than in the other overdominance models.

Species	Author	Number of animals tested
African buffalo	Wenink et al. (1998)	150
Alpine chamois	Schaschl et al. (2012)	364
Pyrenean chamois	Cavallero et al. (2012)	81
Bighorn sheep	Gutierrez-Espeleta et al. (2001)	213
Soay sheep	Paterson et al. (1998)	1209
Scottish Blackface sheep	Schwaiger et al. (1995)	299
Chinese goat breeds	Li et al. (2006)	459
Norwegian Red cattle	Kulberg et al. (2007)	523
Japanese Black cattle	Miyasaka et al. (2011)	338
Iranian Holstein cattle	Nassiry et al. (2005)	250

Table 4.3: Bovidae populations used for calculating measures of allelic diversity in section 4.4.2.6. Only allele frequency data was used from these articles, as sequence data was generally not available.

I compared these data to emergent features of the population generated by the evolutionary simulations. Specifically, I compared expected heterozygosities, the share of the gene pool of the five most frequent alleles and the diversity profile ${}^qD^I$ (Leinster and Cobbold, 2012). A diversity profile is a continuum of diversity measures, which differentially weight the relative abundance of distinct alleles, obtained by varying the viewpoint parameter q from zero to infinity. The diversity profile encompasses a range of standard diversity measures: $q = 0$ gives species richness or allele number; $q = 1$ corresponds to Shannon entropy; $q = 2$ corresponds to the Simpson index or expected homozygosity; and $q = \infty$ corresponds to Berger-Parker or the frequency of the most common allele at a locus. This approach provides a robust means of comparing diversities.

4.4 Results

4.4.1 Model verification

To validate the model, I used analytical results from two simpler models rather than the overdominance models discussed, as no analytical results are available for the latter. One of these simpler models was a neutral model, where there is no difference in allele intrinsic merit and no heterozygote advantage, while the other one was a fully symmetric overdominance (FSOD) model that is symmetric both in the fitness of homozygotes and heterozygotes (i.e. $f_{ij} = d > \tilde{d} = f_{ii} \quad \forall i, j (i \neq j)$). Two quantities were compared, one of which was the expected homozygosity, and the other one was the distribution of allele frequencies (the “allele frequency spectra”).

Table 4.4 compares the share of homozygotes obtained from the simulations for neutral alleles and a FSOD model to analytical values obtained from Kimura and Crow (1964). All results were averaged over k repeats at the end of the simulation span; the number of repeats is specified in table 4.4. Although homozygosity levels in the neutral model were more variable, simulated values were still in good agreement with theoretical expectations. For the SOD model, there was little variation in the amount of homozygosity between the repeats, and thus this model was in much closer agreement with the analytical results of Kimura and Crow (1964).

Neutral model						
population size	10000	25000	100000	250000	1000000	2500000
repeats (k)	50	80	40	20	15	8
homozygosity (theor.)	96.15%	90.91%	71.43%	50.00%	20.00%	9.09%
homozygosity (simu.)	96.21%	92.76%	69.94%	54.48%	20.08%	10.40%
FSOD model						
population size	25000	10000	25000			
selection coefficient (s)	0.05	0.25	0.25			
repeats (k)	20	50	20			
homozygosity (theor.)	4.75%	3.71%	2.27%			
homozygosity (simu.)	4.78%	3.72%	2.26%			

Table 4.4: Comparison of simulated heterozygosities to theoretical values for a neutral model and a FSOD model using the formulae provided by Kimura and Crow (1964).

In addition to the comparison of the expected homozygosities I compared the distributions of allele frequencies (the “allele frequency spectra”, Maruyama and Nei (1981)). This was done by using results from Kimura and Crow (1964) and Maruyama and Nei (1981) in case of the neutral model (figure 4.4), and Kimura and Crow (1964), Yokoyama and Nei (1979) and Nei (1987) for the FSOD model (figure 4.5). The shapes of these frequency spectra obtained from analytical results were in good agreement with frequency spectra obtained from stochastic simulations of the respective models (figures 4.4 and 4.5). For the neutral model, the transition from a U-shaped to a L-shaped distribution was reflected in the simulations (figure 4.4). For the symmetric overdominance model, figures 4.5c and 4.5d further demonstrate that the theoretical frequency spectra still described the allele frequency distribution well even when some amount of variation among the intrinsic merits of alleles is added (i.e. if the FSOD model becomes a SOD model), as long as heterozygote advantage was the dominant force of selection.

These results validated the simulation program for the two simple cases of neutrality and fully symmetric overdominance. The other overdominance scenarios investigated, particularly simple and reinforcing asymmetric overdominance and divergent allele advantage can not strictly speaking be verified by theoretical results. However, the model mechanics are not affected by the choice of overdominance model. The only difference between the overdominance models is the conversion of intrinsic merits of alleles and allele sequence differences (in case of the DAA model) to the fitness of the resulting genotype. I extensively checked these conversions to ensure that the conversion routine reflected the intended behaviour.

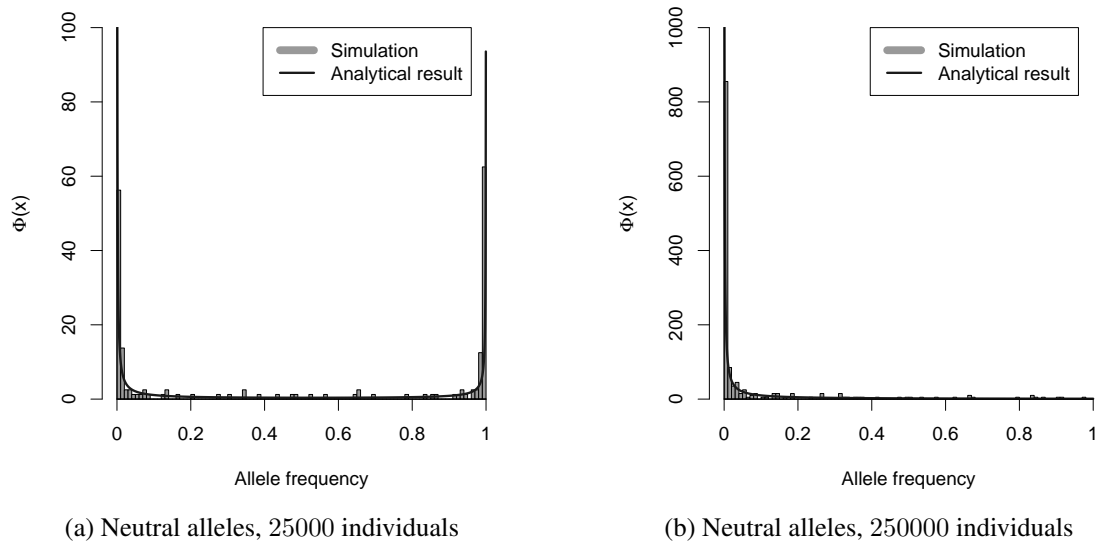
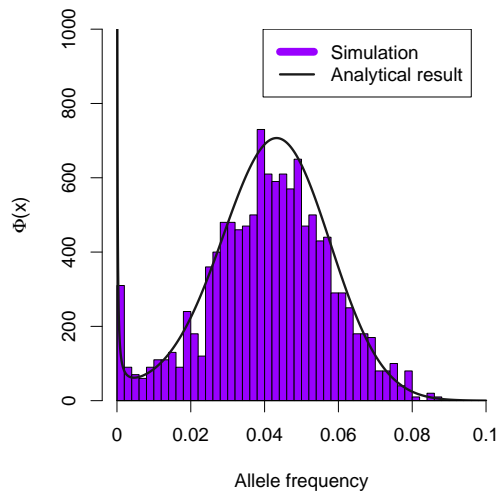
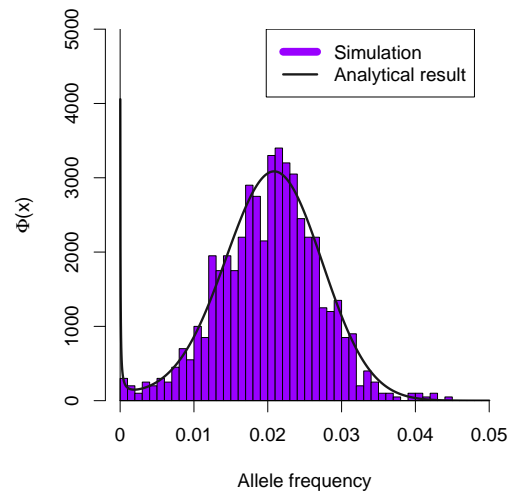
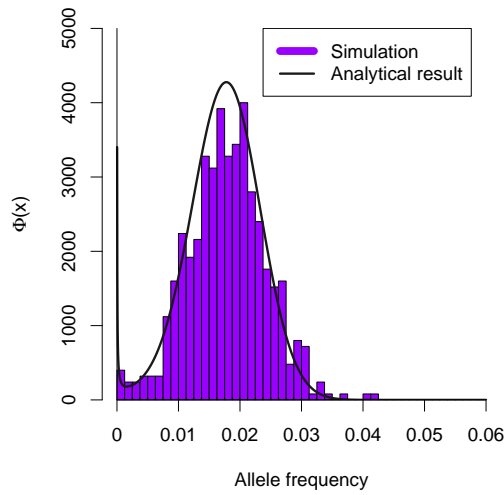
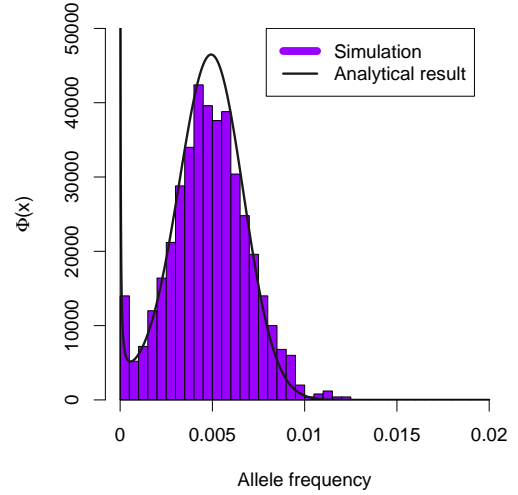


Figure 4.4: Comparison of frequency spectra of neutral alleles resulting from the stochastic simulations to analytical results.

(a) FSOD model, 25000 individuals, $s = 0.05$ (b) FSOD model, 25000 individuals, $s = 0.25$ 

(c) SOD model, setting 2, 25000 individuals



(d) SOD model, setting 2, 250000 individuals

Figure 4.5: Comparison of frequency spectra of the FSOD model and the SOD model of setting 2 to analytical results.

4.4.2 Comparison of Overdominance Models

4.4.2.1 Number of alleles

The number of alleles is one of the defining characteristics of MHC polymorphism, and any theory attempting to explain selection at the MHC is expected to predict large numbers of alleles. However, caution needs to be exercised when drawing conclusions from the sheer number of alleles that a model trying to explain the MHC polymorphism predicts. While the number of alleles is a popular measure of diversity, it is still only a single statistic, and importantly, it does not take allele frequencies into account. I address in detail later the fact that this may result in an overestimation of allelic diversity in some scenarios, if most or nearly all alleles occur at very low frequencies, becoming extinct after existing for only a short time span (these alleles will subsequently be called “transient” alleles, see section 5.3.1.3).

The stochastic simulations of allele evolution showed that the number of alleles in a single, well-mixed population initially increased rapidly (figures 4.6 and B.4) and then entered a relatively stable phase for most of the simulation span of 10 million generations. Alleles were continually created by mutation, but the vast majority were subsequently lost by selection and drift.

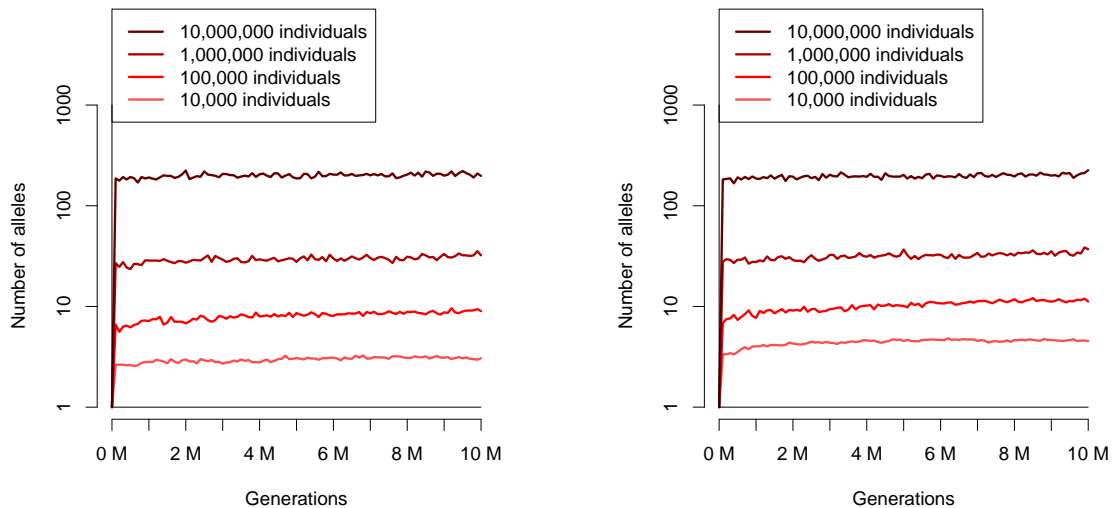


Figure 4.6: Number of alleles over time in the DAA model for settings 1 (left) and 2 (right). Shades of red represent different sized populations. The simulation starts with a single allele and models the action of mutation, selection and drift over millions of generations.

Figures 4.7 and B.5 show the number of alleles that were present at the end of the simulation (after 10 million generations of evolution) for different population sizes. For small population sizes, the stochastic loss of alleles (genetic drift) is the main factor determining allele numbers (Kimura, 1983), whereas the variation in allele intrinsic merit, the amount of

heterozygote advantage and the mutation rate become more influential for larger populations (Li, 1978). However, the different overdominance models showed substantially different responses to an increase in population size.

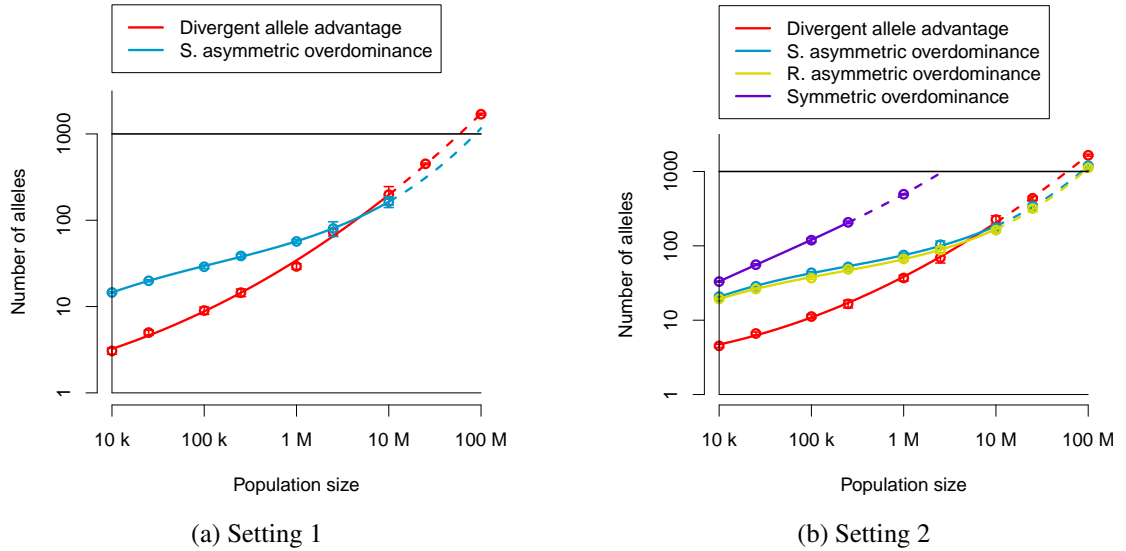


Figure 4.7: Number of alleles vs. population size for settings 1 and 2 and different overdominance models.

The SOD model is known to be capable of maintaining large numbers of alleles (Kimura and Crow, 1964; Maruyama and Nei, 1981; Takahata et al., 1992); this was confirmed by these simulations. Even when allowing for a moderate amount of variation in the fitness of homozygotes, a thousand alleles could be generated for population sizes of only a few million (figure 4.7b).

The two asymmetric overdominance models (SAsOD and RAsOD) had a very similar response in terms of allele numbers as a function of population size. Allele numbers were relatively high for small population sizes, but the rate of increase was only modest when populations became larger.

The DAA model showed relatively low numbers of alleles for small populations, but a stronger increase in allele numbers when populations became larger, compared to the asymmetric overdominance models. The number of alleles maintained under this model exceeded 1000 for population sizes between 50 and 100 million across a broad parameter range (figure 4.8).

These results show that all the overdominance models compared could maintain large numbers of alleles when population sizes were large enough and mutation created new alleles at a constant rate. Allele numbers reached a thousand alleles for population sizes between 1 and 100 million for all the overdominance models, and over a wide range of parameter values. The exact number of alleles present in a population was dependent upon all the pa-

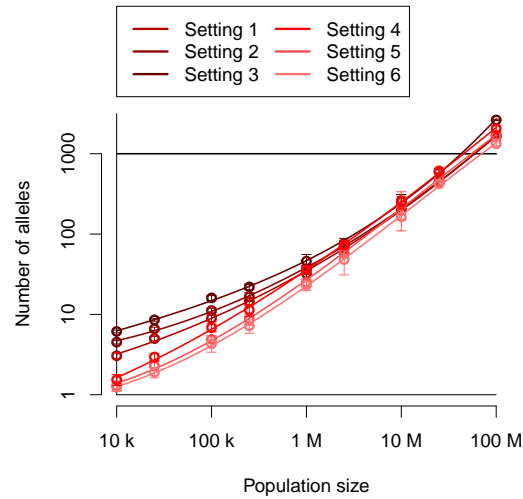


Figure 4.8: Population size and the number of distinct alleles in a population (DAA model). The number of alleles present after 10 million generations varied with choice of σ (variation in intrinsic merit of alleles) and δ (relative heterozygote advantage per amino acid difference). Six different settings ((σ, δ) -pairs, see table 4.2) were examined.

rameters: the population size (m), the mutation rate (μ), the number of variable sites (n) and the relative strengths of selection among alleles (σ) and for heterozygotes (δ). As for the last two parameters, the largest number of alleles occurred when the heterozygote advantage was relatively large compared to the strength of selection among alleles (figure 4.9 shows the response in allele numbers for the DAA model and a population size of 1 million). This is in line with earlier results based on equilibrium analysis (Lewontin et al., 1978).

4.4.2.2 Allele frequency distributions

As the number of alleles maintained did not allow for discrimination between different overdominance models, the next model characteristic examined was the distribution of the frequencies of the final set of alleles. If no selection is present (i.e. if a neutral model is assumed), then a characteristic pattern of allele frequencies, resulting from mutation and random loss of alleles, can be expected. If overdominance is present, however, then this pattern will change. Differences in the frequency distributions of alleles between loci where overdominance is present and neutral loci are often used to detect balancing selection (Ewens-Watterson test, Ewens (1972), Watterson (1978)). Alternatively, characteristics of allele frequency distributions could be used to discriminate between different overdominance models.

In general, very uneven distributions with one dominant allele, as for example in situations with pure dominance, will increase the probability of stochastic loss of alleles for the low-frequency alleles. Models with strongly uneven distributions will consequently be very sensitive to population size, and this overshadows other selection mechanisms that might modify

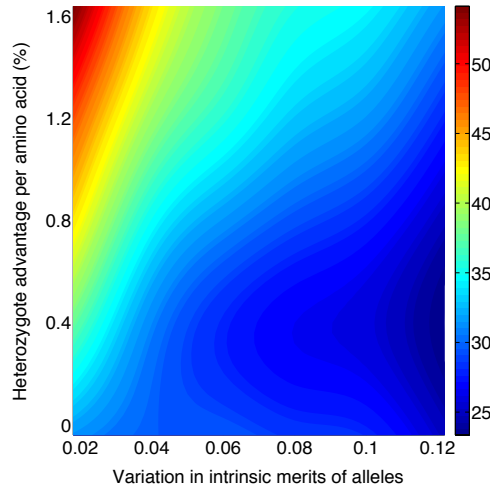


Figure 4.9: Number of alleles for different sets of parameter values in the DAA model. Allele numbers (indicated by colours) depended upon variation in the intrinsic merit of alleles (calculated as the coefficient of variation) and heterozygote advantage (calculated as the gain in fitness of the heterozygote individual per amino acid difference between its alleles). Colours indicate the number of alleles after 10 million generations, averaged over 5 repeats. The population size was 1 million.

the number of alleles maintained in small populations.

One method to compare different overdominance models is to plot the allele frequency distributions for single simulation runs and compare the shapes of these distributions. However, as this method compares results of single simulation runs, it is prone to stochastic fluctuations. A more powerful and robust method of summarising the frequency data in a population is to count the number of representatives of all alleles that belong to a particular frequency class. For example, if three alleles exist in a population with frequencies of 0.95, 0.04 and 0.01, and the allele pool is separated into 10 frequency classes, then the class of the relatively rare alleles with frequencies between 0 and 0.1 has a frequency of $0.04 + 0.01 = 0.05$, and the class of the high-frequent alleles with frequencies between 0.9 and 1 has a frequency of 0.95, while all the other classes in between have a frequency of 0. In general, if j classes are used, then alleles with frequencies below $\frac{100}{j}\%$ form the lowest class, alleles with frequencies between $\frac{100}{j}\%$ and $2 \cdot \frac{100}{j}\%$ the second lowest class and so on, until all alleles are assigned to one of the j possible classes. This classification was made for every simulation run, and the class frequencies were averaged over all runs.

Alleles with similar frequencies were thus accumulated into classes by adding the number of representatives of each allele in a class together. This approach provided signatures of allele frequency distributions and is similar to the “allele frequency spectrum” introduced by Ewens (1972), which ultimately builds on examinations of allele frequency distributions by Wright (1990) and provides the basis for the Ewens-Watterson test of neutrality (Watterson,

1978). The only difference from the method presented here is that classical frequency spectra count alleles but do not weight them according to their frequencies (i.e. number of allele representatives that are present in the population). By accounting for the frequencies of alleles, or, in other words, by counting allele representatives, a spectrum can be obtained that is less sensitive to stochastic fluctuations, particularly if results of multiple runs are averaged. It furthermore emerged that the shape of these weighted frequency spectra was similar for all population sizes examined. This approach therefore provided a robust tool for evaluating the evenness of allele frequencies.

Figure 4.10 shows the weighted frequency spectra of the overdominance models as well as for a neutral model and a model of pure dominance for population sizes of 100000.

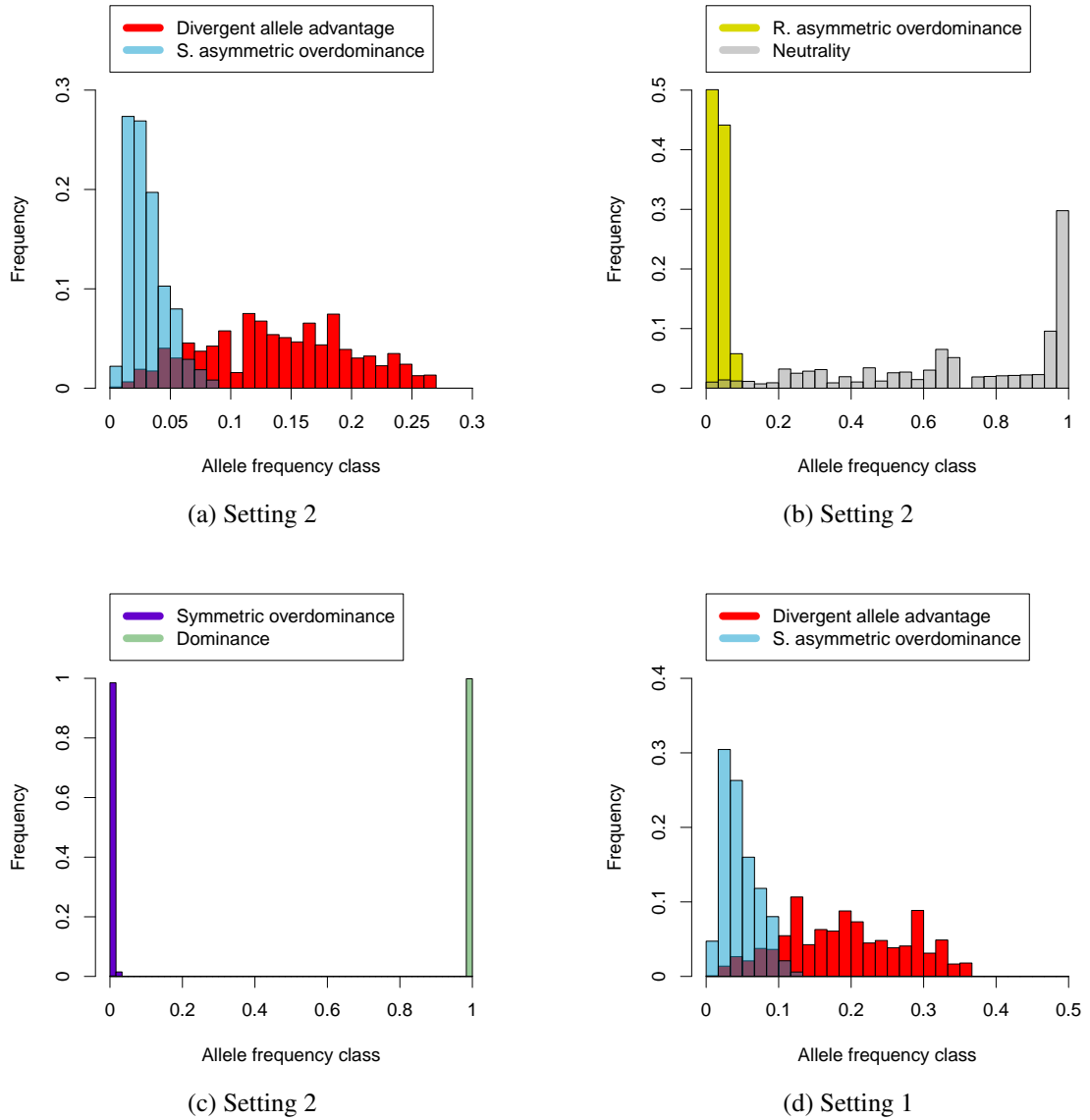


Figure 4.10: Weighted frequency spectra for different models and settings 1 and 2. The population size was always 100000.

The different models compared in figure 4.10 showed a very diverse range of weighted allele frequency profiles. Under pure dominance, the dominant (fittest) allele had a frequency close to 100% (figure 4.10c). If neutrality of alleles was assumed, then the allele frequencies became more even, but still only a few alleles were relatively abundant, while the majority of alleles occurred at low frequencies (figure 4.10b). In the DAA model the allele frequencies were again more even than in the neutral model, but not as even as allele frequencies were in the simple and reinforcing asymmetric overdominance models (figures 4.10a, 4.10b, 4.10d, B.7, B.8, B.9 and B.10). The SOD model was ultimately the model with the most even allele frequencies (most alleles fell into the lowest frequency classes, figures 4.10c and B.6).

While observed allele frequencies for some MHC loci are known to be more even than under a neutral model (Hedrick and Thomson, 1983), the evenness exhibited by the SOD model clearly exceeded reported values in most cases (see the data presented in table 4.3), except in situations where the amount of heterozygote advantage was exceptionally low compared to the variation in allele intrinsic merit. These observations already cast doubt on the SOD model as a fundamental force that drives allelic diversity at the MHC. Furthermore, there is some recent experimental evidence that heterozygote advantage is acting asymmetrically, i.e. some heterozygote genotypes are fitter than others (Bronson et al., 2013), in contrast to the assumptions of the SOD model. Therefore, I will limit most of the subsequent comparisons to the DAA and asymmetric overdominance models.

The exact shape of the weighted frequency spectra, or how even the allele frequencies of each model are, depended on the setting examined. For example, allele frequencies were less even under symmetric overdominance if the amount of heterozygote advantage was very low compared to the variation in allele intrinsic merit, as in setting 5 (figure B.6). However, the qualitative shape of the weighted frequency spectrum of each model and the order of the models in terms of allele evenness, reflected by these frequency spectra, were similar for the majority of settings and population sizes (figures B.7, B.8 and B.9). Only for low population sizes (less than 100000 individuals) in settings with a low amount of heterozygote advantage (settings 4, 5 and 6), the DAA model regularly resulted in monomorphic populations, therefore yielding a qualitatively different weighted frequency spectrum (figure B.10a).

The obtained weighted frequency spectra could in principle be used to compare simulation results to experimental allele frequency data, and thus distinguish between candidate models. However, as the choice of number and range of frequency classes is somewhat arbitrary, I postpone the comparison of allelic diversity between models and data to sections 4.4.2.4, 4.4.2.5 and 4.4.2.6, where easily comparable diversity measures were extracted from allele frequency data.

4.4.2.3 Persisting variation in allele intrinsic merit

Although substantial experimental evidence for heterozygote advantage exists, its relevance for maintaining MHC diversity has been questioned (reviewed in Spurgin and Richardson (2010)). One of the reasons for the reluctance to accept heterozygote advantage as a fundamental driver for MHC diversity is that at equilibrium, the large number of alleles observed can only be maintained if the intrinsic merits of the alleles are in a narrow range, i.e. if there is only low, possibly unrealistically low (De Boer et al., 2004), persisting diversity in allele intrinsic merit. For that reason, I compared models with reference to different measures of allele intrinsic merit variation, and examined the extent to which allele intrinsic merit variation can be maintained over evolutionary time scales.

Maintaining large numbers of alleles did not require markedly similar intrinsic merits in the DAA model, in contrast to the asymmetric overdominance models examined (figures 4.11 and B.11). There was a much wider distribution of intrinsic merits for alleles under the DAA model than under the asymmetric overdominance models.

This can also be demonstrated by comparing the standard deviations of the intrinsic merit of all alleles present in the gene pool at the end of the simulation period. Figures 4.12 and B.19 show the standard deviation in allele intrinsic merit, weighted by the allele frequencies, under different overdominance models.

For a population size of 10 million, the total number of alleles in the DAA model was higher than in the asymmetric overdominance models for all settings except setting 3 (figures 4.7 and B.5), yet the DAA model allowed for a much higher variation in allele intrinsic merit. This contradicts a statement derived from equilibrium analysis that “one obtains a high degree of polymorphism if, and only if, the fitness contributions of MHC alleles are very similar” (De Boer et al., 2004), and is a first indication that allelic diversity is maintained to a much higher degree in the DAA model than in the traditional models of asymmetric overdominance.

In the DAA model, fitness contributions of alleles (i.e. the allele intrinsic merits) can be much more spread out while still allowing for a large number of alleles in a realistic non-equilibrium situation (figures 4.11, B.11, 4.12 and B.19). This is a clear indication that adding the sequence-dependent fitness component had a stabilising effect on those alleles with lower intrinsic merit whose average sequence difference from the other alleles in the gene pool was large enough.

Consequently, alleles with relatively low intrinsic merit occurred at substantial frequencies (figures 4.13, B.20 and B.21) and persisted for long periods (tens of thousands of generations, figures 4.14, B.25 and B.26) provided that they were sufficiently different from the rest of the gene pool. This is in contrast to traditional asymmetric overdominance models, where

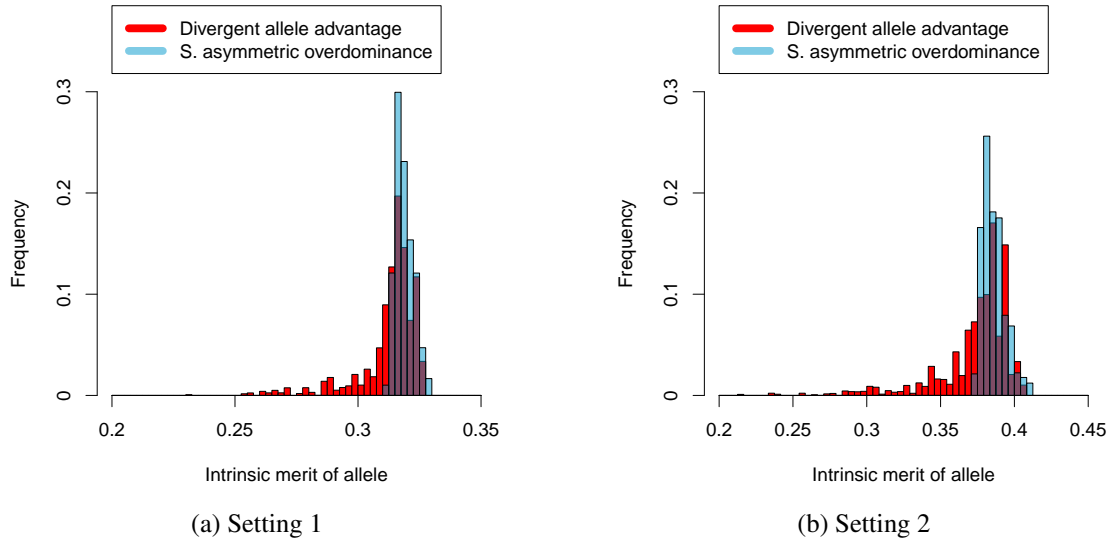


Figure 4.11: Greater variation in allele intrinsic merit under DAA compared to the SAsOD model for setting 1 (left) and setting 2 (right). The frequency distribution is the average over 20 repeats after 10 million generations in the DAA (red bars) and the SAsOD (blue bars) models. Bars refer to the proportion of allele representatives in each interval. The DAA model shows a pronounced “tail” of alleles with low intrinsic merit. The population size was 100000.

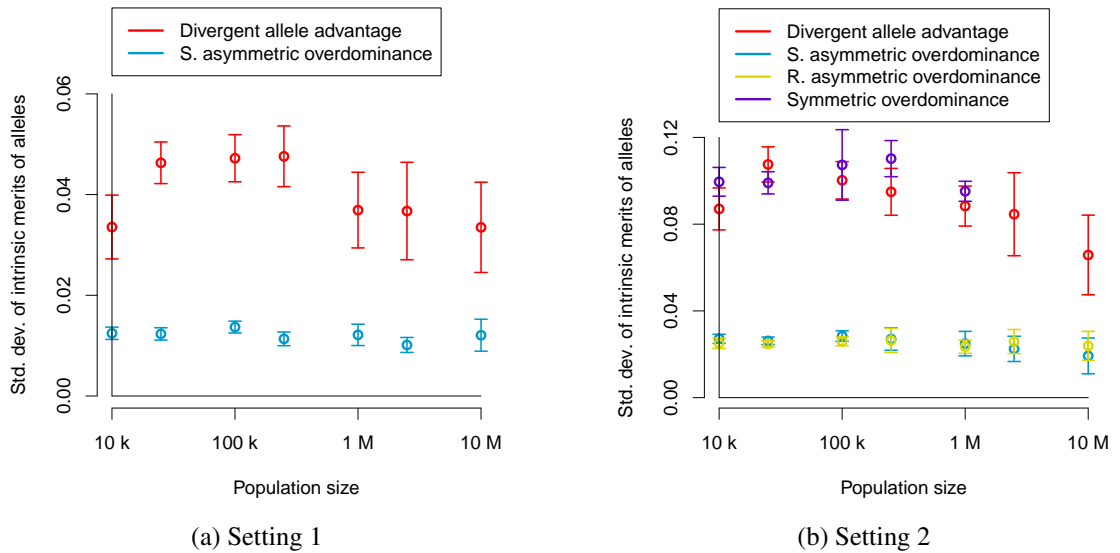


Figure 4.12: Variation in allele intrinsic merit, measured as the standard deviation of the intrinsic merits of the final set of alleles. The alleles were weighted by their proportions. The DAA model is associated with approximately a 3x – 4x larger standard deviation than the asymmetric overdominance models for most settings and population sizes explored.

these alleles were very infrequent (figures B.22, B.23 and B.24).

Different values for the intrinsic merit thresholds w_t that separate alleles with low intrinsic merit from the other alleles in the gene pool were chosen for the six different settings. This was done because differences in the two parameters defining the settings (σ and δ) resulted in allele intrinsic merit distributions with different means and variances for all overdominance models (see figures 4.11 and B.11), and therefore the thresholds were adapted accordingly. Allele intrinsic merit thresholds are discussed in more detail in section 5.3.1.

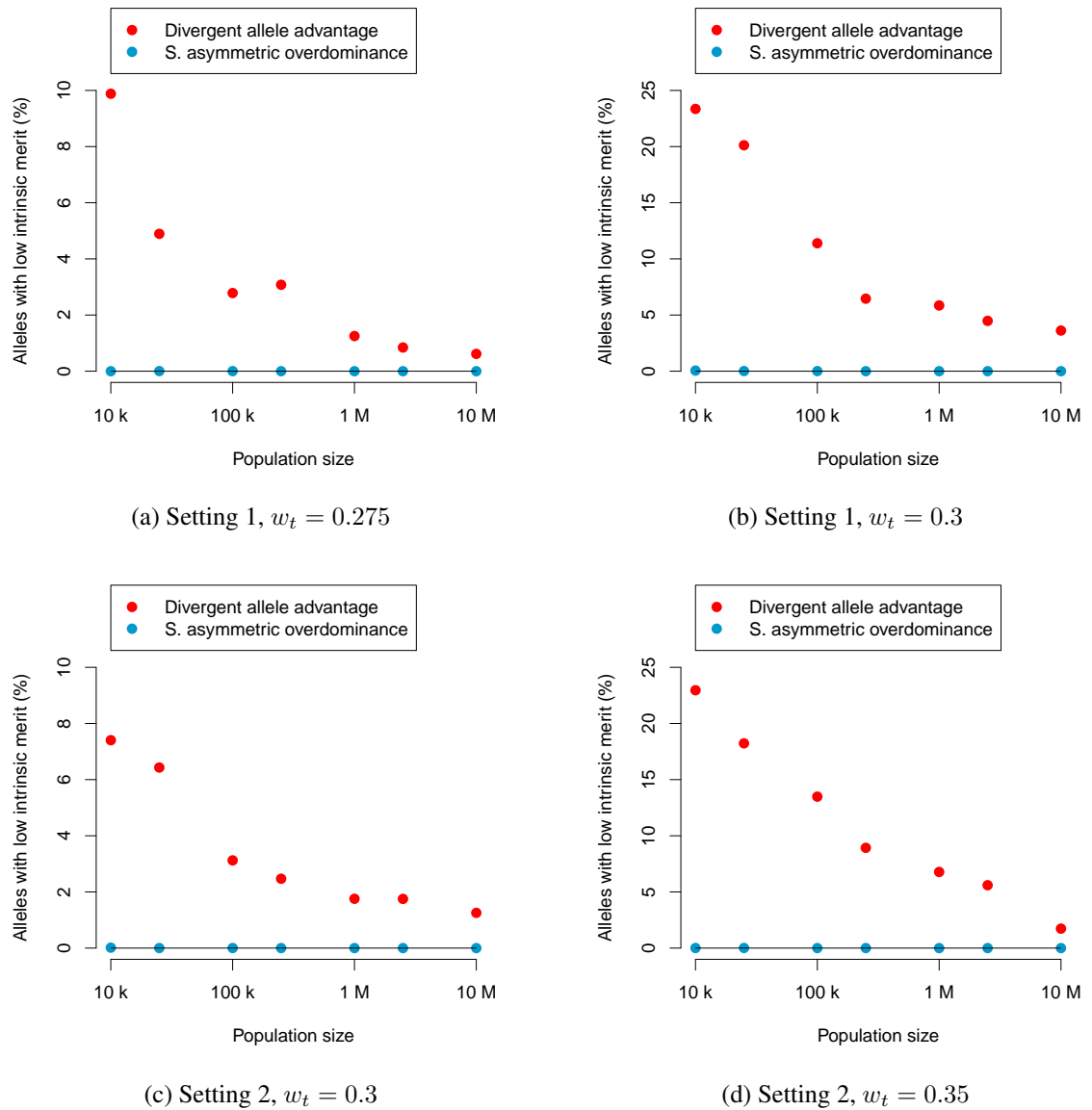


Figure 4.13: Percentage of alleles below two different thresholds of allele intrinsic merit (w_t) for the DAA and SAsOD models and settings 1 and 2.

The much greater variation in intrinsic merit of alleles under DAA compared to traditional asymmetric overdominance models was thus a consequence of the longevity of these alleles, as under both asymmetric overdominance models, alleles with low intrinsic merit only

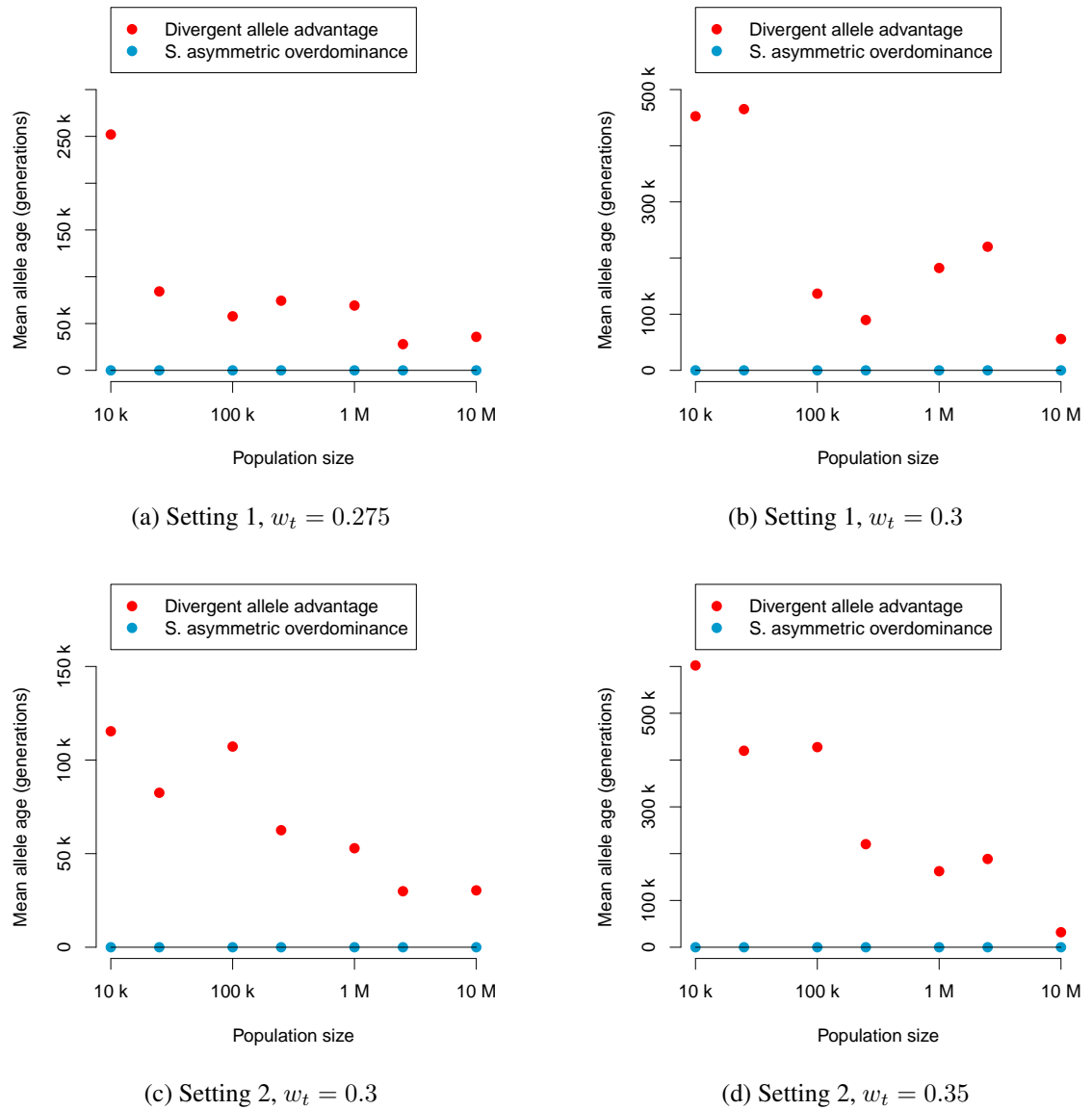


Figure 4.14: Mean age of alleles with low intrinsic merit, using two different intrinsic merit thresholds (w_t), for the DAA and SAsOD models and settings 1 and 2.

existed for short time spans (less than about 150 generations for all settings examined, figures B.27, B.28 and B.29). These alleles formed poor homozygotes, made heterozygotes of below average fitness and were eliminated by natural selection relatively rapidly.

4.4.2.4 Diversity expressed as heterozygosity

Allele numbers, the shapes of allele frequency distributions and the variation in allele intrinsic merit were representations of diversity used to discriminate between models in sections 4.4.2.1, 4.4.2.2 and 4.4.2.3. While there is a large number of diversity measures in ecology and genetics (Hill, 1973; Leinster and Cobbold, 2012), one of the more widely used measures in genetics is the expected frequency of heterozygotes h_e , which can be calculated directly from allele frequency data, assuming that n alleles with frequencies p_i ($i = 1, \dots, n$) are present in the population (equation 4.11).

$$h_e = 1 - \sum_{i=1}^n p_i^2 \quad (4.11)$$

The expected heterozygosity calculated from simulation results was compared against expected heterozygosity calculated from allele frequency data of the DRB locus from ten studies of bovidae populations (five wild populations and five breeds, see section 4.3.7 and table 4.3).

The expected frequency of heterozygotes in these populations of wild and domestic species of bovidae fell between 70% and 90% (figures 4.15 and B.30), indicating that 10%–30% of the population were homozygous at a single MHC locus. This is consistent with the range predicted by the DAA model for population sizes between 10000 and 10 million; in contrast, a much lower frequency of homozygotes was predicted by the traditional asymmetric overdominance models (figure 4.15). This difference in the frequency of homozygotes was repeated across the parameter settings with the exception of those settings that represented exceptionally low levels of heterozygote advantage (figure B.30). The SOD model underestimated homozygosity levels even more (figures 4.15d and B.51a). It must be pointed out that the population sizes used for simulations with the SOD model for calculating diversity measures ranged from 10000 to 250000 only, as only for these population sizes the simulations were performed for the full 10 million generations. This only marginally affected the results, however, as the mismatch between data and model output even increased for higher population sizes.

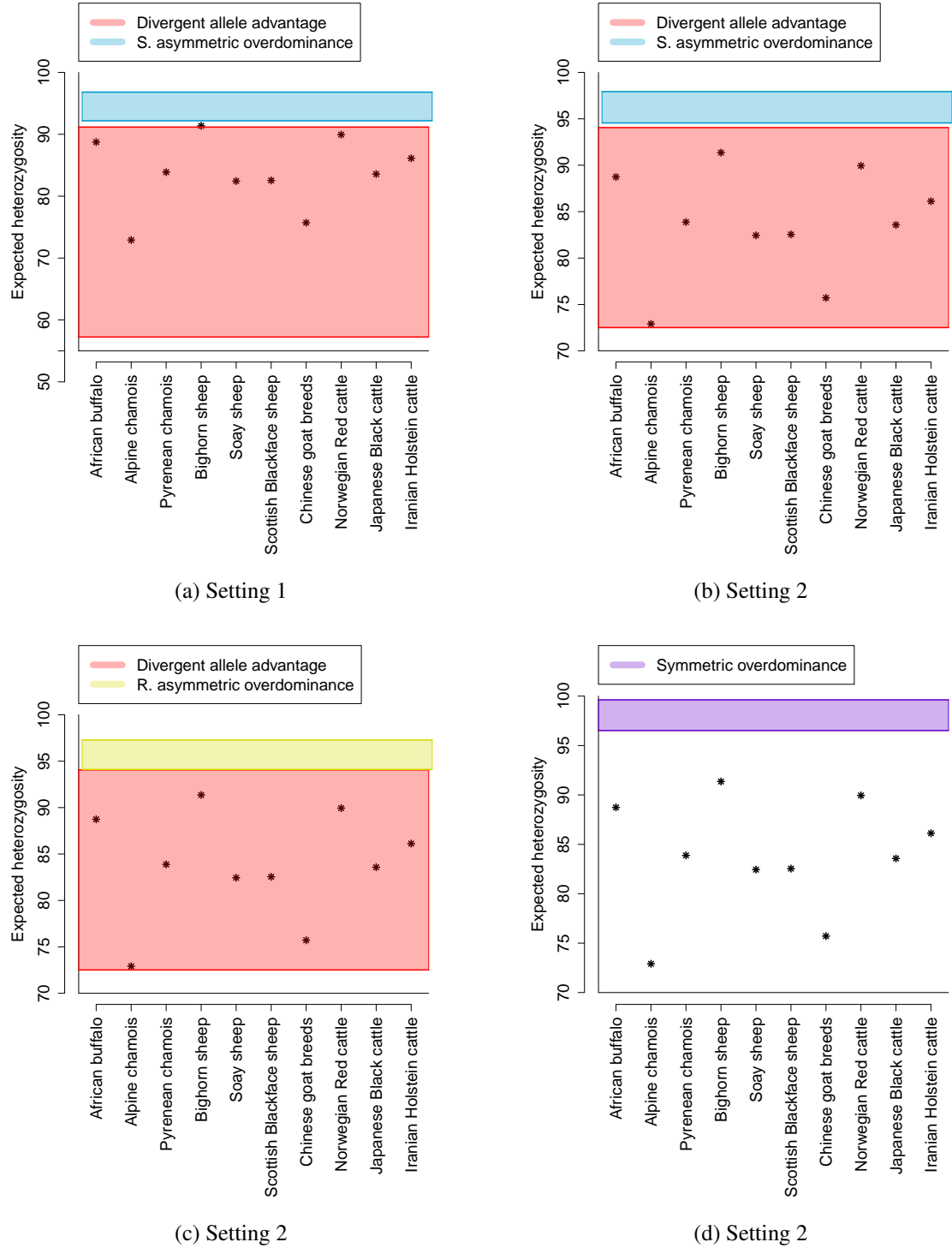


Figure 4.15: Comparison to expected heterozygosities calculated from published allele frequencies for different overdominance models and settings 1 and 2. The expected heterozygosity (in %), represented by stars, is calculated from observed allele frequencies of wild and domesticated populations of bovidae (table 4.3), while model predictions are shown as coloured bands, spanning population sizes from 10000 to 10 million, respectively. A total of 800 alleles were drawn randomly from the final gene pool of each simulation. This represented the average sample size in the observed data of approximately 400 individuals.

4.4.2.5 Diversity expressed as gene pool share of the most frequent alleles

The expected heterozygosity shaped up as a useful measure to compare the evenness of allele frequencies (section 4.4.2.4), and thus allowed for a comparison of model outputs to observed data. Alternatively, one could exclusively consider the frequencies of the k most frequent alleles, with $k \in \mathbb{N}$. Such a diversity measure puts all weight on the more frequent alleles, and completely disregards the alleles with lower frequencies.

Figure 4.16 compares the (cumulative) share of the most frequent alleles of samples from natural populations to simulation results of different overdominance models for $k = 5$ and settings 1 and 2. The examined populations of bovidae species had a share of the 5 most frequent alleles ranging between 50% and nearly 100% (figures 4.16 and B.31). This range was in line with the predictions from the DAA model for settings 1 and 2, whereas other overdominance models failed to predict the observed gene pool share of the 5 most frequent alleles for these settings (figure 4.16).

The same pattern emerged for parameter settings 3 and 4, where the DAA model predicted observed values well while the predictions of the SAsOD model were generally outside the observed cumulative share of the 5 most frequent alleles (figures B.31a and B.31b). As was the case for the heterozygosity, the situation was reversed for settings 5 and 6, which refer to exceptionally low levels of heterozygote advantage. For setting 5, the SAsOD model gave better predictions than the DAA model (figure B.31c), while for setting 6 both of these overdominance models failed to predict observed values (figure B.31d).

4.4.2.6 Diversity expressed as diversity profiles

Both the number of alleles and the heterozygosity level are measures of diversity of the gene pool that were utilised to draw comparisons to observed values. It has been shown that both of these measures, as well as many other measures of diversity, can usefully be included into a continuous family of diversity measures ${}^qD^Z$ (Leinster and Cobbold, 2012) graphically represented by a diversity profile. The “similarity matrix” Z carries information about how different the alleles are. When there is no knowledge of allele similarity, then the identity matrix I can be used instead of the similarity matrix Z , and the diversity profile gives information about the number of alleles and the evenness of allele proportions.

Figures 4.17 and B.32 show a plot of the diversity profile (disregarding allele similarity, $Z = I$) for different overdominance models and a population size of 100000 for all settings, while figures 4.18, B.33, B.34, B.35, B.36 and B.37 show the DAA model and the SAsOD models for population sizes of 10000, 100000, 1 million and 10 million, respectively. For all these plots, the coloured bands indicate the 95% confidence intervals resulting from variation in the simulation repeats.

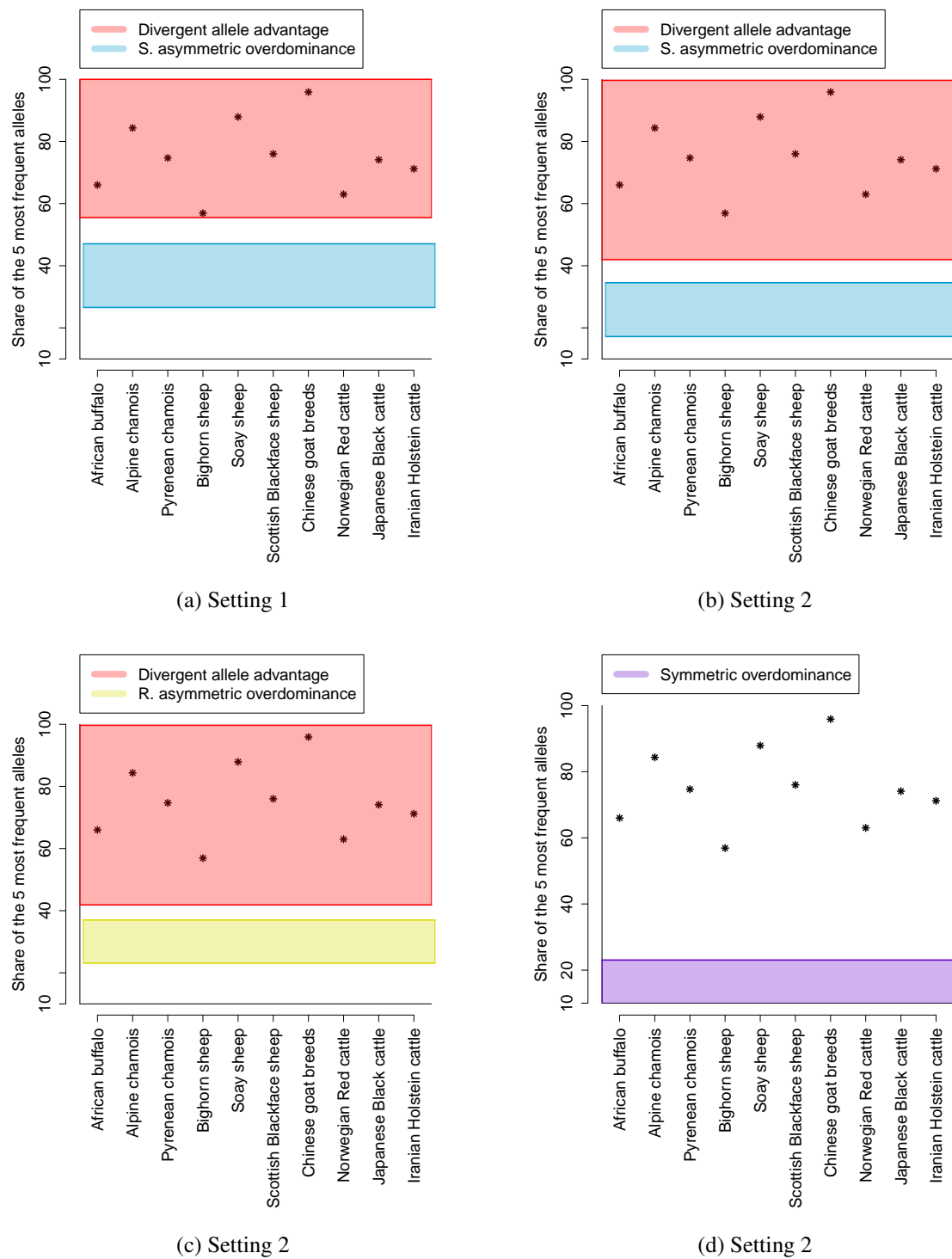


Figure 4.16: Comparison to the share of the 5 most frequent alleles calculated from published allele frequencies for different overdominance models and settings 1 and 2. The share of the 5 most frequent alleles (in %), represented by stars, is calculated from observed allele frequencies of wild and domesticated populations of bovidae (table 4.3), while model predictions are shown as coloured bands, spanning population sizes from 10000 to 10 million, respectively. A total of 800 alleles were drawn randomly from the final gene pool of each simulation. This represented the average sample size in the observed data of approximately 400 individuals.

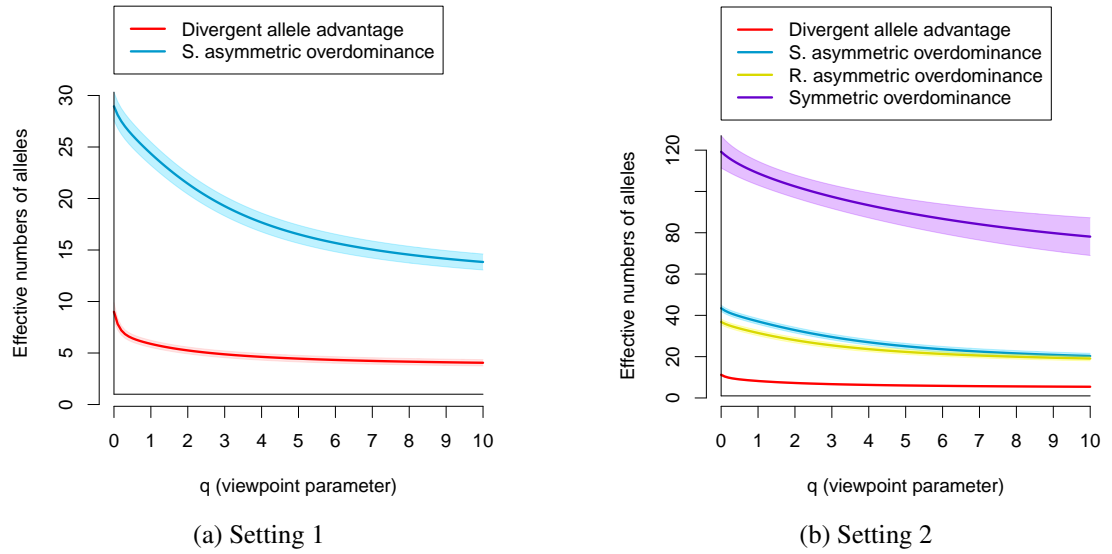


Figure 4.17: Diversity profiles (naive diversity) for different overdominance models and settings 1 and 2. The population size was 100000.

As the diversity measure for $q = 0$ corresponds to allele numbers, the SAsOD model was more “diverse” (i.e. had a higher number of alleles) than the DAA model for lower population sizes, whereas the opposite was true for high population sizes (figures 4.18d, B.33d, B.34d, B.35d, B.36d and B.37d show that for all settings except setting 3, the DAA model exhibited a higher diversity for q close to and including 0 for a population size of 10 million). However, as q increases, allele frequencies played an increasingly important role, and given its more even frequencies, the SAsOD model was more “diverse” for sufficiently large q . This is in line with the earlier observation that the percentage of heterozygotes ($q = 2$) was lower in the DAA model than in the SAsOD model for all population sizes examined (see section 4.4.2.4).

While diversity profiles can usefully be applied to examine evenness of alleles, one of their main advantages can only be exploited when similarity between alleles is taken into account. The similarity between any two alleles is combined in a “similarity matrix” \mathbf{Z} . The elements of the similarity matrix \mathbf{Z} may, for example, carry information about the difference in allele intrinsic merit of two alleles i and j . An element of \mathbf{Z} is therefore calculated by subtracting the modulus of the difference in intrinsic merit between alleles i and j from the maximum possible intrinsic merit difference (0.5), and dividing by the maximum possible intrinsic merit difference (see equation 3.2).

Alternatively, the entries of the similarity matrix \mathbf{Z} could be based on the number of differences between the two encoded amino acid sequences on the considered locus according to equation 4.12:

$$\mathbf{Z}_{ij} = \frac{n - k_{ij}}{n} \quad (4.12)$$

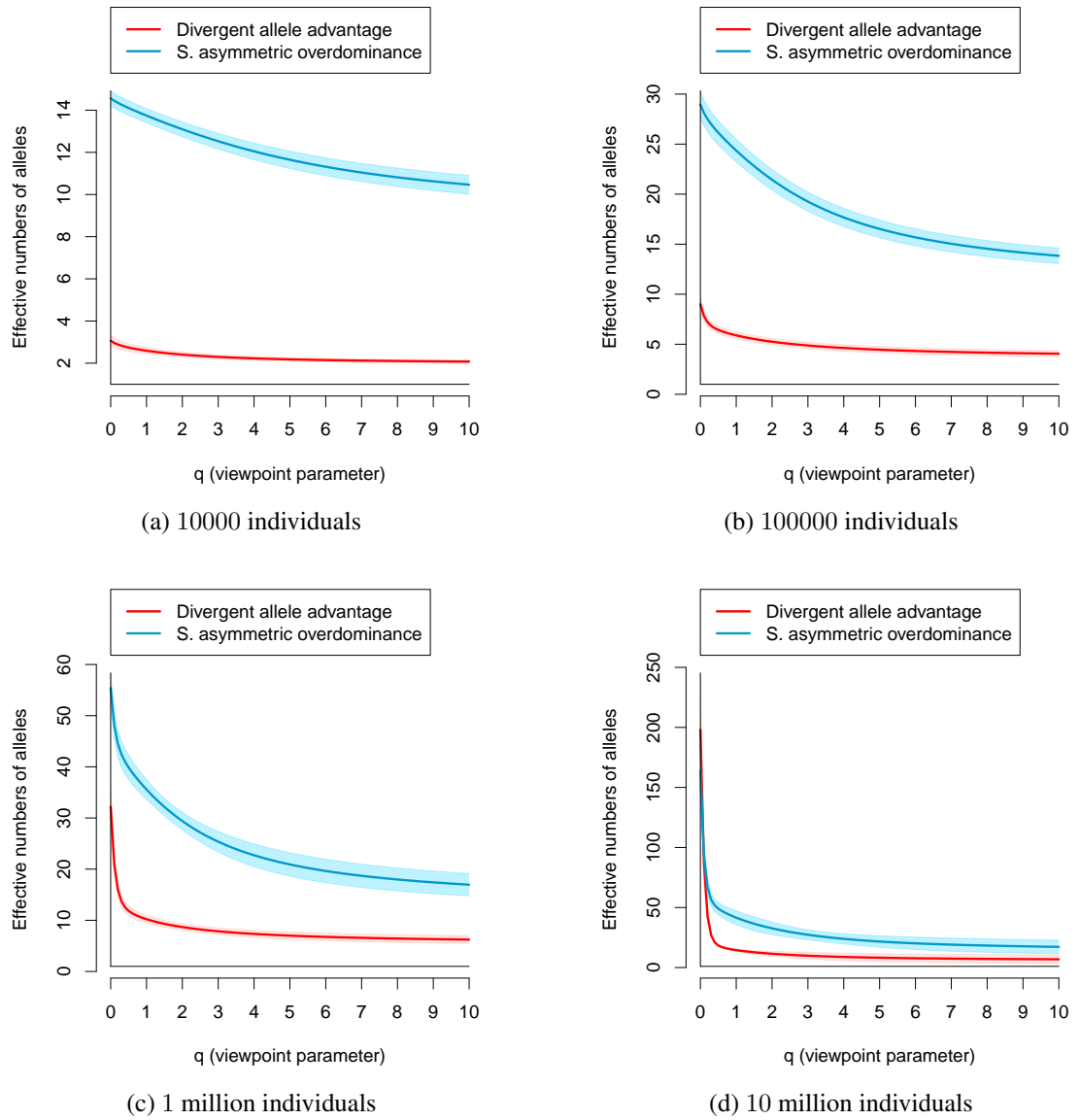


Figure 4.18: Diversity profiles (naive diversity) for the DAA and the SAsOD models, setting 1 and population sizes of 10000, 100000, 1 million and 10 million.

Here, n denotes the number of variable positions on the protein chain encoded by an allele (with 82 as default value), while k_{ij} is the actual number of amino acid differences between two alleles A_i and A_j .

If the identity matrix I is substituted by a similarity matrix Z based on allele intrinsic merit difference, sets of alleles that show greater intrinsic merit differences between the alleles can be expected to be more diverse, given a similar frequency evenness of the models. Even though the SAsOD model had much more even allele frequencies (figure 4.17), the allelic diversity when taking allele similarity into account was clearly higher in the DAA model (figure 4.19). This was true for all values of q , for all settings $((\sigma, \Delta)$ -pairs) examined (figures 4.19 and B.38) and all population sizes (figures 4.20, B.39, B.40, B.41, B.42 and B.43), as long as the gene pool was polymorphic. Only when small populations were considered, and the amount of heterozygote advantage was low, was the allelic diversity that takes allele similarity into account lower for the DAA model. This applied to settings 4 and 5 for population sizes below 100000 and setting 6 for population sizes up to 100000 (figures B.38d, B.41a, B.42a, B.43a and B.43b). These exceptions are not surprising, as the previously mentioned scenarios led to a monomorphic gene pool in a considerable number of simulation repeats.

These findings demonstrate that the results presented in the previous section, namely that the DAA model exhibits a considerably larger intrinsic merit variation in the final gene pool, is highly robust.

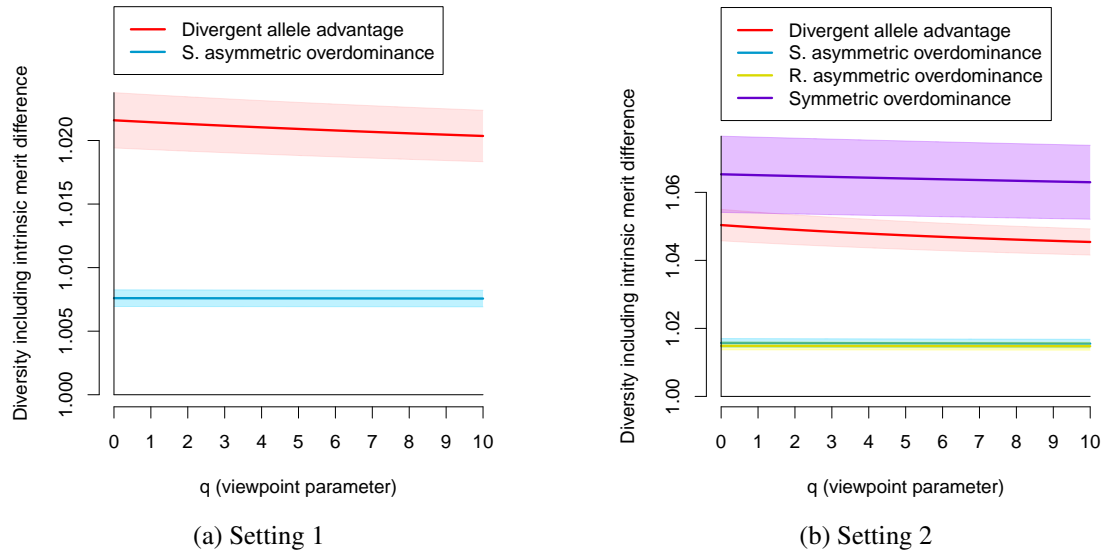


Figure 4.19: Diversity profiles when accounting for intrinsic merit differences between alleles for different overdominance models and settings 1 and 2. The population size was 100000.

What is more, exactly the same effect can be seen when replacing intrinsic merit differences between alleles (equation 3.2) by differences in the amino acid sequences encoded by the

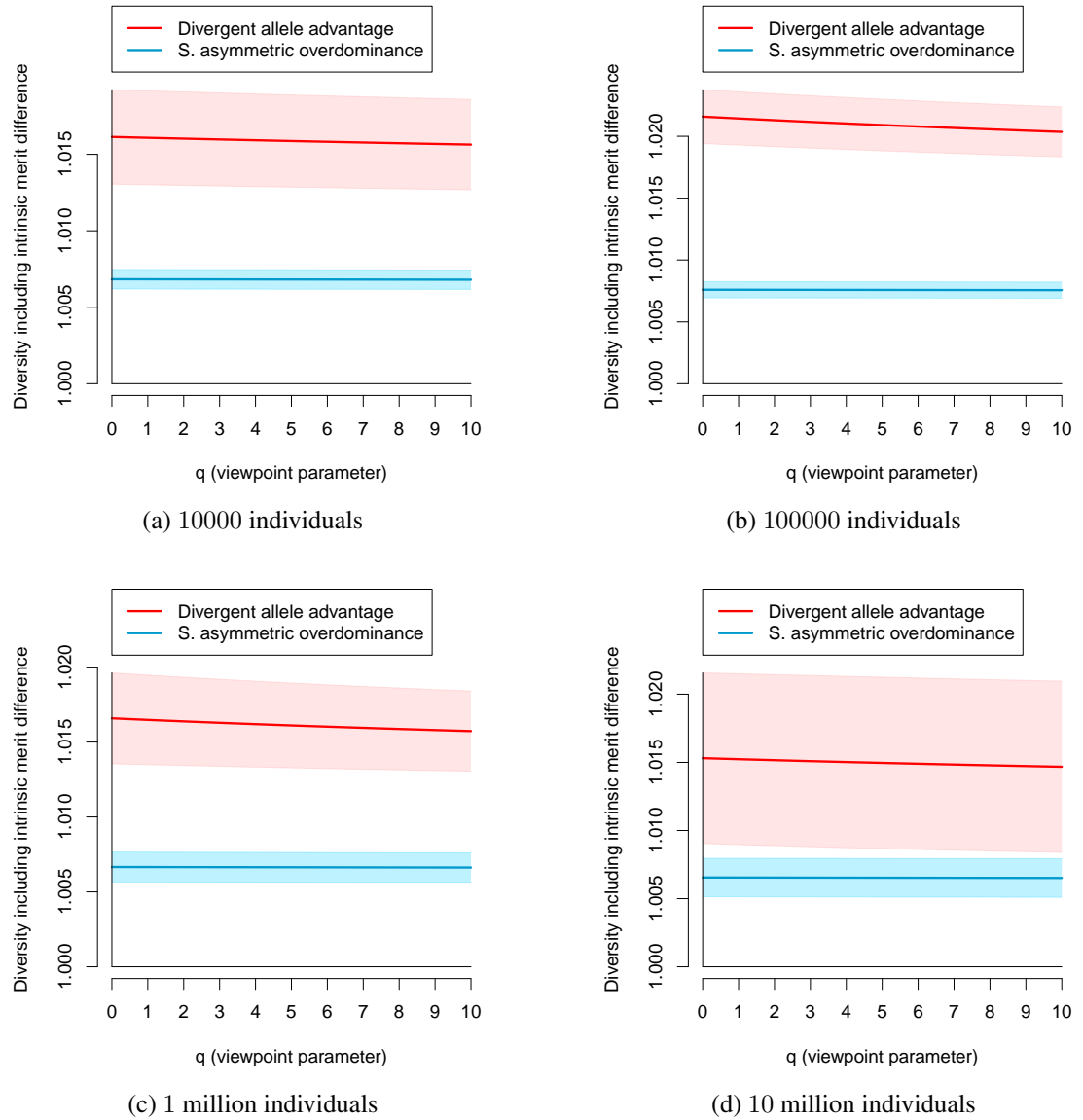


Figure 4.20: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and the SAsOD models, setting 1 and population sizes of 10000, 100000, 1 million and 10 million.

alleles (equation 4.12) as the diversity measure (figures 4.21, 4.22, B.44, B.45, B.46, B.47, B.48 and B.49).

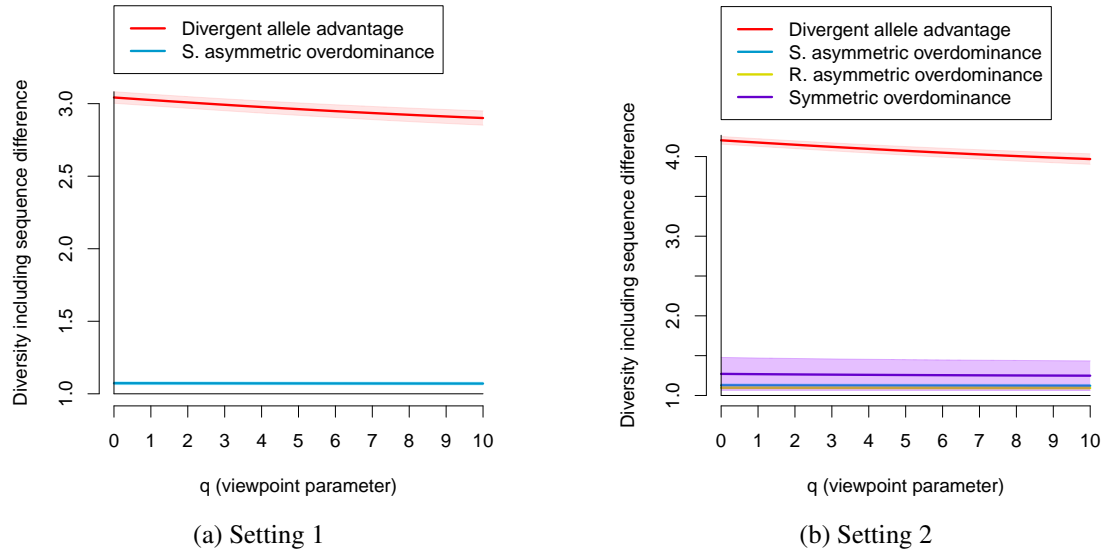


Figure 4.21: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for different overdominance models and settings 1 and 2. The population size was 100000.

While the previous observations confirmed the fact that the DAA model can maintain a larger variation of allele intrinsic merit and more diverse allele sequences in the gene pool of a population, the question still remains whether it also matches observed allele frequencies when comparing diversity at multiple scales. The diversity profiles calculated with samples from natural populations (figures 4.23 and B.50) demonstrated that most measures of diversity for most populations fell within the range predicted by the DAA model. A simple explanation beside parameter uncertainty for the relatively small number of observations falling outside the predicted boundaries is that the intensity of selection may differ among populations. Models with different intensities of selection produced different diversity profiles. Notably, models with a very low ratio of heterozygote advantage to variation in allele intrinsic merit (i.e. settings 5 and 6) did not produce diversity profiles similar to those observed in natural populations (figures B.50c and B.50d).

The SOD model in turn substantially overestimated the evenness of allele frequencies (expressed as naive diversity) for all values of q , even for setting 5, which is a setting with a low ratio of heterozygote advantage to allele intrinsic merit (figure B.51).

4.4.2.7 Trans-species polymorphism

Trans-species polymorphism is another feature exhibited by MHC genes. Some alleles are old enough to be shared by different species. The simplest explanation for this observation

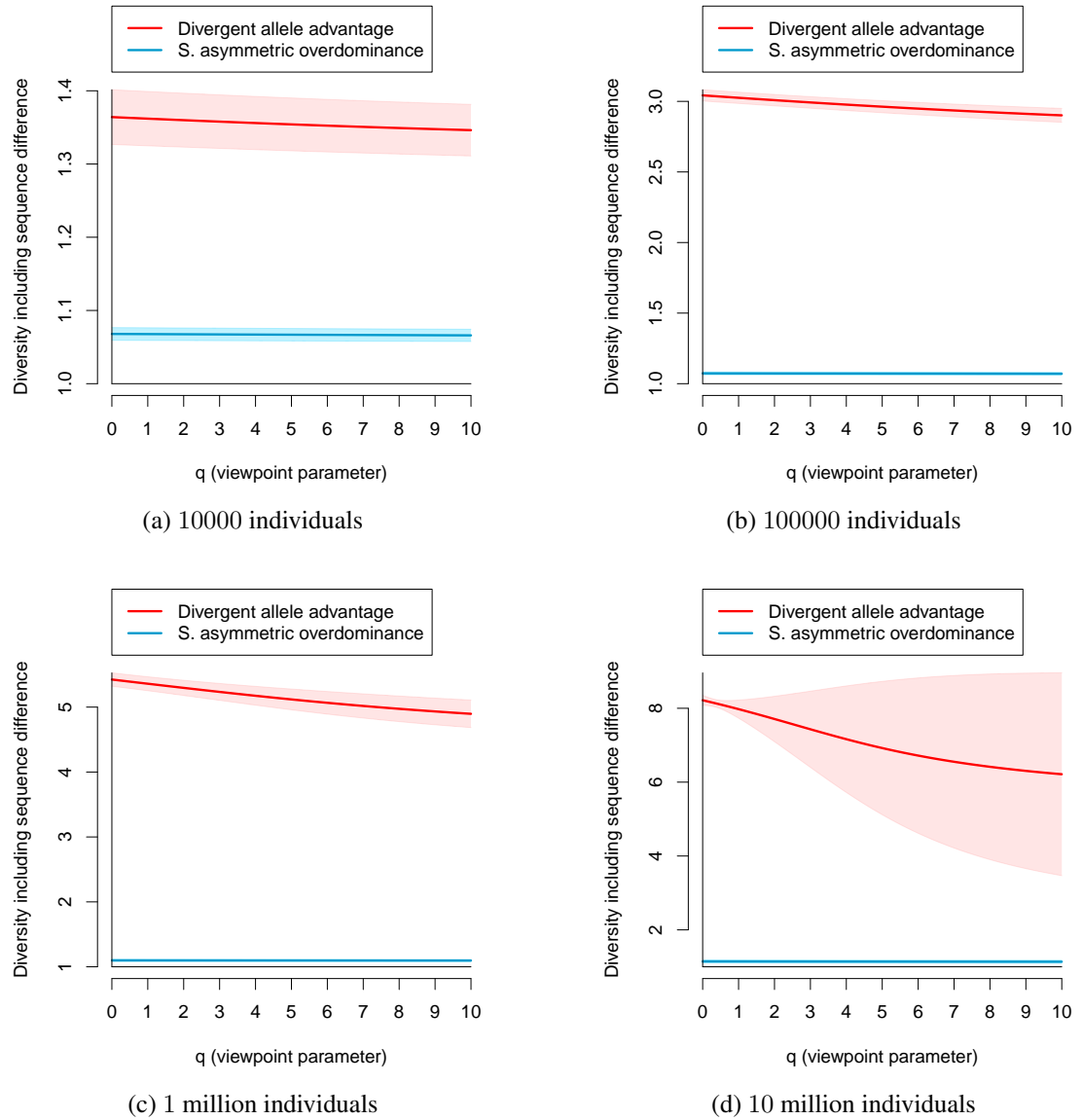


Figure 4.22: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and the SAsOD models, setting 1 and population sizes of 10000, 100000, 1 million and 10 million.

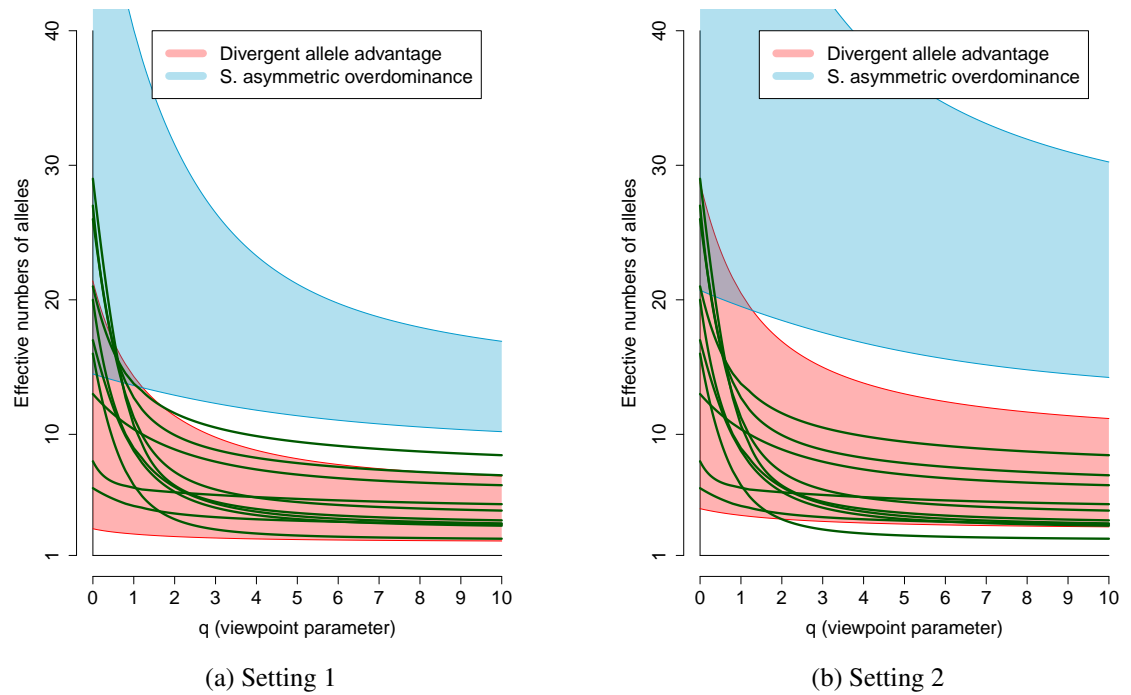


Figure 4.23: Diversity expressed as profiles, showing the diversity profiles calculated from observed allele frequencies of wild and domesticated populations of bovidae (green curves) and as predicted by alternative models of heterozygote advantage. Coloured bands show the range of diversity profiles predicted by the DAA model (red band) and the SAsOD model (blue band), spanning population sizes from 10000 to 10 million, respectively. 800 alleles were drawn from each simulated population to reflect the sampling of the observed populations.

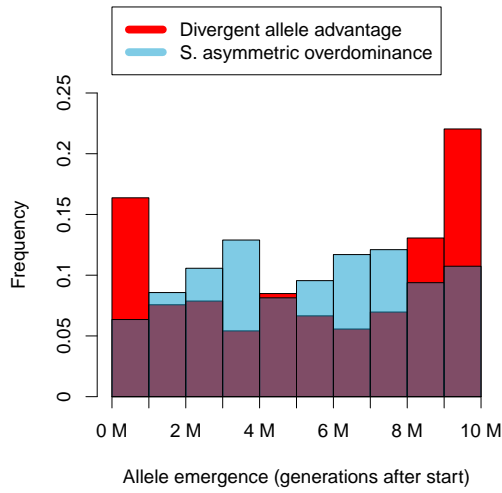
is that these shared alleles arose prior to speciation. In the DAA model, the final gene pool at the end of the simulation exhibited a mix of alleles that emerged early and alleles with more recent emergence times (figure 4.24); this was true for all parameter settings and population sizes (figures B.52, B.53, B.54, B.55, B.56 and B.57). The lengthy existence times of some alleles shown in these figures is consistent with trans-species polymorphism. Alleles with high intrinsic merit that arose early in the simulation were hard to displace while among the more recent alleles were a substantial number of alleles with low intrinsic merit; these alleles had not yet been eliminated by natural selection. The distribution of the emergence times of alleles in the DAA model is therefore typically U-shaped. The size of the arms representing new and old alleles varied with the size of the population. Small populations had a relatively large number of new alleles while larger populations had a relatively large number of old alleles (figures B.52, B.53, B.54, B.55, B.56 and B.57).

Experimental data shows that alleles currently present in the bovidae species can be attributed to different periods of evolution, as ancient trans-species allelic lineages have been identified in a number of MHC genes (Ballingall et al., 2010; Sena et al., 2011). This means that a viable model should allow some of the alleles that were created in the early stages of the 20 – 40 million years of evolution to persist to the end of the simulation. In fact, it even should allow some alleles generated throughout the different periods of evolution to be present in the final gene pool. While such a result was found for the DAA model and the asymmetric overdominance models (figures 4.24, B.52, B.53, B.54, B.55, B.56, B.57 and B.58), the SOD model struggled to explain the existence of old alleles (figures 4.24 and B.59) unless the amount of heterozygote advantage was very low compared to the variation in allele intrinsic merit. This provides yet another reason for challenging this model as a plausible representation for the generation and maintenance of MHC diversity.

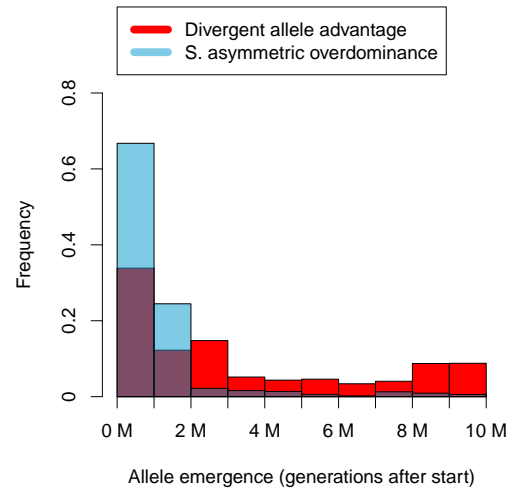
4.4.3 Sensitivity towards mutation rate and number of amino acids per allele

The key parameters of the stochastic simulations, outlined in section 4.3.5, were the population size (m), the variation in intrinsic merit of alleles (σ), the relative heterozygote advantage per amino acid difference (δ) as well as the mutation rate (μ) and the number of amino acids per allele (n). The first three parameters are associated with a large amount of uncertainty in real-world populations. Therefore the sensitivity of all results towards variation in m , σ and δ was extensively explored by simulating populations ranging from 10000 to 10 million, and checking the results for different settings, i.e. different (σ , δ)-pairs.

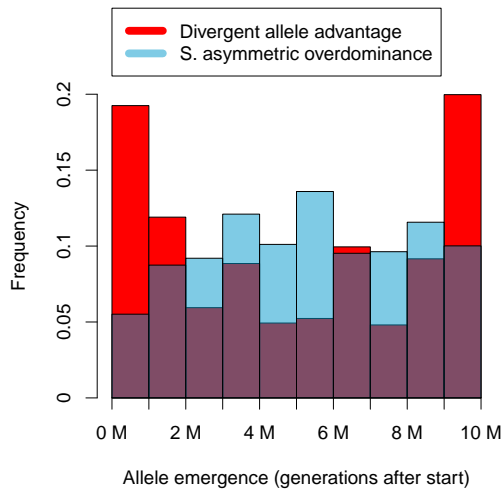
The other two parameters, mutation rate and number of amino acids per allele, are associated with better estimates. While some authors argue for a similar mutation rate for all mammals



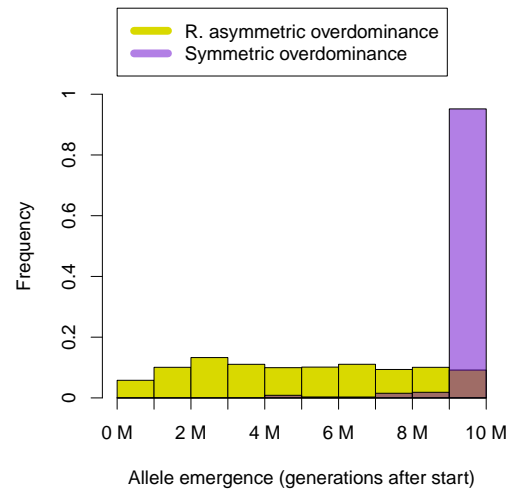
(a) Setting 1, 100000 individuals



(b) Setting 1, 10 million individuals



(c) Setting 2, 100000 individuals



(d) Setting 2, 100000 individuals

Figure 4.24: Distribution of emergence times of alleles for different overdominance models and settings 1 and 2.

(Kumar and Subramanian, 2002) and regions of the genome, others found evidence for differences between species (Wu and Li, 1985), even between closely related ones (Ohashi et al., 2006). Furthermore, there is indication for differences in mutation rates for different regions of the genome (Nachman and Crowell, 2000; Roach et al., 2010) or the sex of an individual (Nachman and Crowell, 2000).

Likewise, the number of amino acids per allele could be interpreted in different ways. First, n could be viewed as the exact number of amino acids of a particular allele, no matter whether these sites are in a peptide binding region or not. This viewpoint gave rise to the default value for n when considering the DRB1 locus of bovids. Second, n could be viewed as the approximate number of amino acids in the peptide binding regions of that locus. Such an approach was taken by e.g. Satta (1997) and Lenz (2011), and results in a lower estimate for n . Third, as some closely linked MHC loci, the DRB1, DQA1 and DQB1 loci of bovids among them, nearly exclusively appear as haplotypes (Andersson et al., 1988), n could also be viewed as the number of amino acids in the peptide binding regions of the entire haplotype, or even as the total number of amino acids in the haplotype. This could increase the value of n above the default.

These considerations imply a necessity to check the validity of the main results when μ and n take on different values. I varied both μ and n by halving and doubling the standard value listed in table 4.1 to explore trends in results when decreasing or increasing μ and n , respectively.

The fitness benefit per amino acid difference (σ) was adjusted when n was varied such that the product $\sigma \cdot n$ remained constant. The sensitivities of the results when varying μ and n were explored for setting 1 and population sizes of 10000, 25000, 100000, 250000, 1 million, 2.5 million and 10 million.

Expected heterozygosities and diversity profiles for setting 1 and different mutation rates and numbers of variable sites were calculated from the frequencies of alleles in the final gene pools of the simulation runs in the same way as in sections 4.4.2.4 and 4.4.2.6.

Figures 4.25 and B.60 show the sensitivity of the number of distinct alleles for different population sizes. When the mutation rate was varied (figures 4.25a and B.60a), it emerged that a higher mutation rate (μ) increased the number of alleles present after 10 million generations for all population sizes. Reducing the number of variable sites (figures 4.25b and B.60b), however, led to a higher number of alleles at small population sizes, whereas increasing it led to a higher number of alleles at large population sizes.

Figures B.61 and B.62 show the sensitivity of the range of allele intrinsic merits and the weighted standard deviation in allele intrinsic merit for different population sizes, respectively. Varying the mutation rate affected the range of allele intrinsic merits only marginally (figures B.61a and B.61c) and slightly decreased the weighted standard deviation in allele

intrinsic merit when the mutation rate increased (figures B.62a and B.62c). Varying the number of amino acids per allele again had little effect on the range of allele intrinsic merits (figures B.61b and B.61d) and only decreased the weighted standard deviation in allele intrinsic merit for lower population sizes when the number of variable sites was increased (figures B.62b and B.62d), while larger population sizes were hardly affected by a change in the number of variable sites.

Figures 4.26, B.63 and B.64 show the frequency distribution of allele intrinsic merits for setting 1 and different population sizes, while figures B.65 and B.66 show the same information for setting 2. Again, the sensitivity towards mutation rate and number of variable sites was explored. The frequency distribution of allele intrinsic merits was more affected when the number of variable sites was varied (figures 4.26b, B.63b, B.64b and B.64d), with a larger number of variable sites resulting in a more stretched-out distribution, but was only marginally affected when the mutation rate was varied (figures 4.26a, B.63a, B.64a and B.64c).

The frequency distribution of emergence times of alleles for different population sizes, mutation rates and number of variable sites is shown in figures 4.27, B.67 and B.68 (setting 1) as well as figures B.69 and B.70 (setting 2). The qualitative shape of the distribution of the emergence times of alleles was only marginally affected by both the variation in mutation rate and the number of variable sites for all population sizes.

Finally, model outputs related to gene pool diversity for different mutation rates and number of variable sites were compared against experimental data from table 4.3. Figure 4.28 compares the expected heterozygosity and figure 4.29 the diversity profiles. The coloured bands represent the range of population sizes from 10000 to 10 million.

The boundaries of the expected heterozygosity window as well as diversity of the sampled 800 alleles for all other values of the viewpoint parameter q were only slightly shifted when the mutation rate was increased, but were decreased when the number of variable sites was increased.

Taking all these observations together, I conclude that although dependencies on both the mutation rate and the number of variable sites exist, the main results are valid over a considerable range for these two parameters. Diversity measures are generally more affected by a variation in the number of variable sites than by a variation in the mutation rate; however, only a substantial increase in the number of variable sites leads to a marked decrease in diversity as given by the diversity profile of the allele sample after 10 million generations.

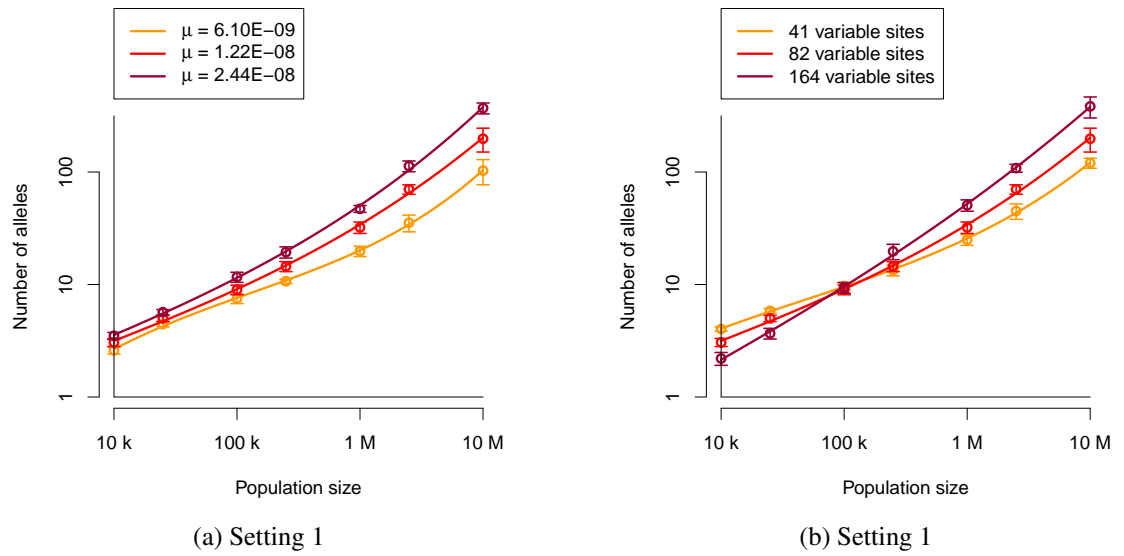


Figure 4.25: Number of alleles vs. population size for setting 1 and different mutation rates (left) and number of variable sites (right).

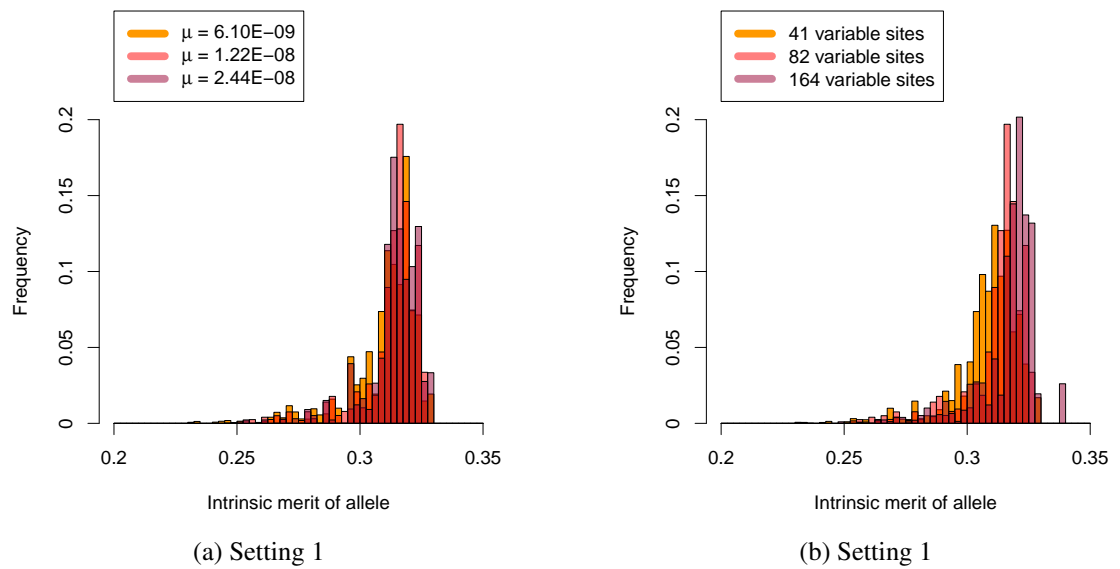


Figure 4.26: Frequency distribution of allele intrinsic merits for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 100000.

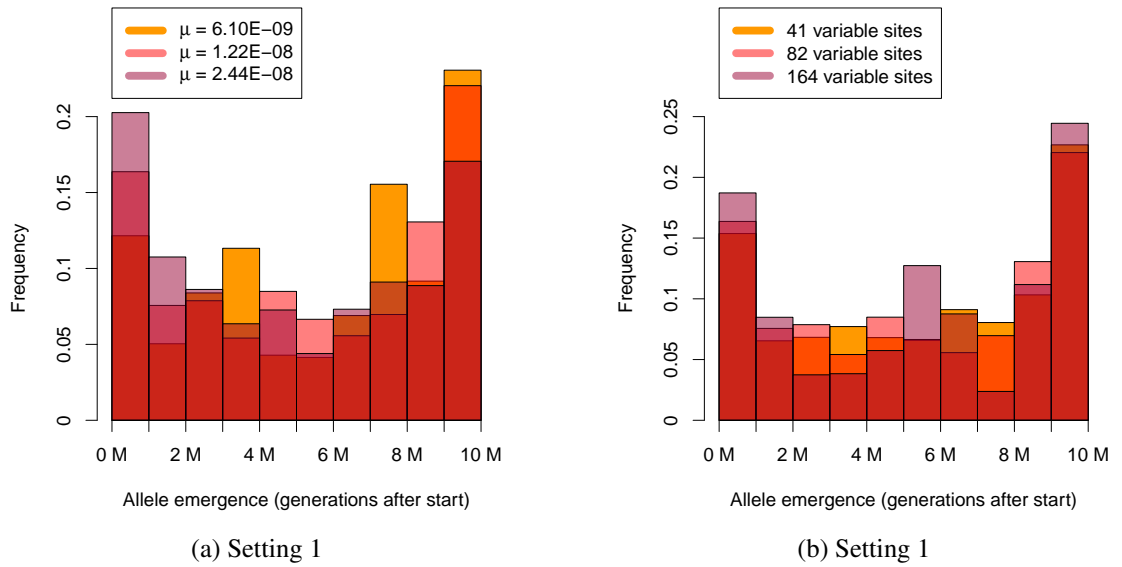


Figure 4.27: Frequency distribution of emergence times of alleles for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 100000.

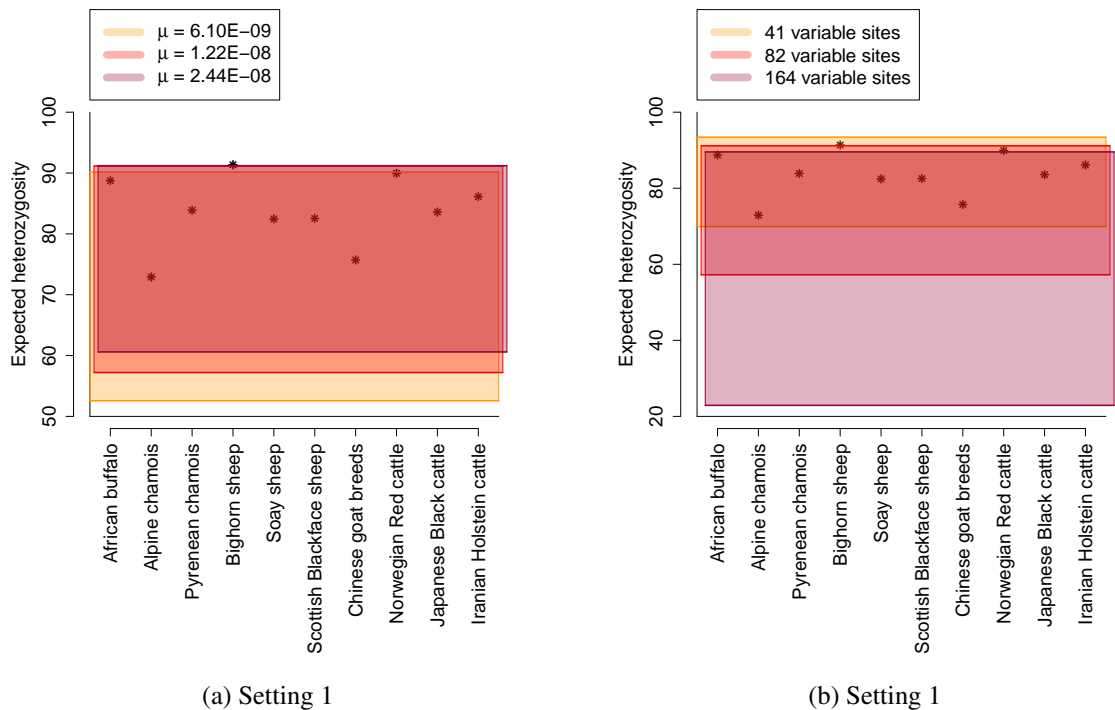


Figure 4.28: Expected heterozygosity for setting 1 and different mutation rates (left) and number of variable sites (right). Model outputs (coloured bands) were compared against experimental data (see table 4.3). The bands represent the range of population sizes from 10000 to 10 million.

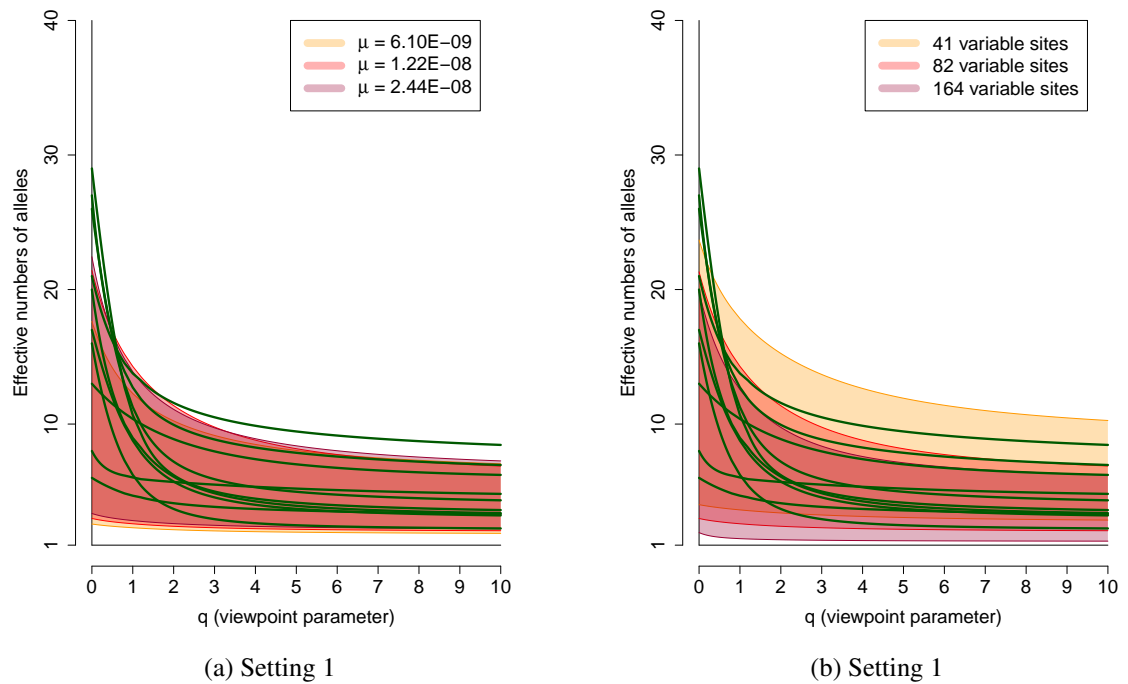


Figure 4.29: Diversity profiles for setting 1 and different mutation rates (left) and number of variable sites (right). Model outputs (coloured bands) were compared against experimental data (see table 4.3). The bands represent the range of population sizes from 10000 to 10 million.

4.5 Discussion

I have used a unique feature of MHC evolution in the bovidae to recreate the generation of MHC diversity from a single ancestral sequence. The traditional asymmetric forms of heterozygote advantage did not reproduce the observed features of MHC diversity; the models only maintained a large number of alleles when the alleles were very similar in their fitness contributions, and the frequency of homozygotes was lower than observed in sampled populations. In contrast, the DAA form of heterozygote advantage produced results that were similar to all the major observed features of MHC diversity, including the large number of alleles, the substantial proportion of homozygotes, the variation in fitness contributions of alleles and the similarities among alleles from different species. The mechanistic biological underpinning of the DAA model, and its consistency with the observations while being parsimonious when compared to alternative explanations of MHC diversity make divergent allele advantage a promising candidate for being a fundamental driver of MHC polymorphism.

The importance of heterozygote advantage as the major selective force maintaining MHC polymorphism has been undervalued, even though there is a clear rationale for overdominant selection based on the broader immune surveillance of heterozygotes compared to homozygotes (Hughes and Yeager, 1998). Other forms of selection, such as frequency-dependent

selection driven by host-parasite coevolution (Borghans et al., 2004), mate choice (Potts et al., 1994) or Associative Balancing Complex evolution (van Oosterhout, 2009) may modify the patterns of MHC diversity (Hedrick, 2002). However, my results demonstrate that it is not necessary to assume an essential role for these mechanisms. In some cases these mechanisms may act in concert with DAA; for example host-parasite coevolution could modify the fitness landscape in which DAA acts.

Traditionally, low MHC diversity in certain species such as bison, moose and Syrian hamsters (Yuhki and O'Brien, 1990; Mikko et al., 1999) has been explained by demographic processes such as a recent bottleneck in the number of breeding individuals (Menotti-Raymond and O'Brien, 1993). Demographic processes may be the fundamental cause of low diversity but low diversity could also arise by selection. MHC diversity is influenced by the relative intensity of selection for heterozygotes compared to selection among alleles (figure 4.9). This ratio will depend upon the severity of the parasite challenge and the sequences of the MHC alleles. Relatively weak selection for heterozygotes would produce low MHC diversity.

4.6 Conclusion

In conclusion, it was demonstrated that divergent allele advantage has the potential to be a fundamental driver of MHC diversity. This special form of overdominance explains observations related to MHC alleles, as large allele numbers, a medium level of evenness of allele frequencies, variation in fitness among individuals, and trans-species polymorphism substantially better than traditional overdominance models. This also implies that it is not necessary to invoke other mechanisms than overdominance to explain these features of MHC alleles – while these mechanisms may contribute, their contribution is not essential.

Chapter 5

MHC Diversity and Disease Susceptibility

5.1 Introduction

Many loci in the vertebrate genome can be characterised by a traditional view of Darwinian selection at the molecular level, resulting in the survival of the allele that conveys the highest level of fitness in the particular environment, while other alleles decrease in frequency and eventually vanish. New alleles are created by mutation, but only occasionally increase to higher frequencies when they confer an above average fitness to the organism.

This traditional view does not apply when intra- or inter-locus interactions exist. One of the most important intra-locus interactions is heterozygote advantage, also referred to as overdominance. With this selective force present, a much more diverse pool of alleles can exist in a population, both in an equilibrium scenario (chapters 2 and 3) and in a scenario where alleles are subject to the evolutionary forces of mutation, selection and drift (chapter 4). This is a common feature of overdominance (Kimura and Crow, 1964; Takahata and Nei, 1990), and was true for all models of heterozygote advantage examined.

Chapter 4 demonstrated that a model based on the divergent allele advantage (DAA) hypothesis can explain the major features of MHC variation. In this chapter, some of the implications this new understanding of MHC functioning has will be illustrated.

While overdominance is a powerful mechanism in maintaining a diverse pool of alleles (Maruyama and Nei, 1981), chapter 4 has demonstrated that splitting the fitness contributions of alleles into an overdominance-independent part, referred to as intrinsic merit, and an additional part that adds the fitness benefit of heterozygotes (based on the underlying overdominance model) results in qualitatively different distributions of allele intrinsic merits for different overdominance models considered. In traditional asymmetric overdominance mod-

els only alleles in a narrow band of intrinsic merits can persist, and recent mutants with an intrinsic merit below that band only exist at low frequencies and for very short time spans. In contrast, the DAA model allows for much greater variation in allele intrinsic merit, and features a pronounced tail of viable alleles with low intrinsic merit that can and do exist for long time spans (see section 4.4.2.3).

These qualitatively different intrinsic merit spectra can be used to discriminate between overdominance models (as was done in chapter 4); in this chapter, however, the characteristics of the alleles with low intrinsic merit will be explored in-depth (section 5.3.1). The results of this analysis will form the basis of a number of important implications that this new understanding of MHC functioning has.

5.2 Materials and Methods

In this chapter, the same overdominance models were used as in chapter 4. Stochastic simulations were performed in the same way and with the same constraints as in chapter 4, and the same parameters and parameter ranges were used.

5.3 Results

5.3.1 Classifying alleles with low intrinsic merit

To understand why the DAA model behaves differently to the asymmetric overdominance models, it is necessary to take into account the key feature characterising divergent allele advantage: the difference in amino acid sequences between any two alleles. This turns out to be of great help in understanding why alleles with low intrinsic merit are so much more long-lived and abundant in the DAA model than in traditional asymmetric overdominance models, and this is the ultimate reason for the much larger variation in the intrinsic merit of persisting alleles in the DAA model (see section 4.4.2.3).

5.3.1.1 Allele intrinsic merit thresholds

As the first aim of this chapter is to classify alleles with low intrinsic merit, i.e. alleles in the tail of the intrinsic merit distributions (see figures 4.11 and B.11), the percentage of allele representatives with a low intrinsic merit was determined. To investigate how stable the obtained results are, this was done for different thresholds for classifying allele intrinsic merit as “low”.

Figures 4.13, B.20 and B.21 show the share of allele representatives for all settings and two different intrinsic merit threshold values. While hardly any representatives of alleles with low intrinsic merit existed in the asymmetric overdominance models, the proportion of alleles with low intrinsic merit in the DAA model was comparatively high, even for the stricter (i.e. lower) threshold value.

Clearly the question remains whether the applied thresholds were suitable for identifying groups of alleles with substantially different behaviour. The aim of applying these thresholds is to separate alleles that are stable in a sense that they persist over substantial periods of time from alleles that vanish quickly after they were created by mutation. Ideally, this partitioning of the set of alleles would not be very sensitive to the exact value of the threshold, and the results would remain the same, qualitatively, if the threshold is varied in a certain range.

An answer to the question posed above can be found by examining the distribution of allele intrinsic merits within a population. Figures 4.11 and B.11 show the distribution of allele intrinsic merits for all settings and a population size of 100000. All settings gave a similar shape for this distribution. The SAsOD model exhibited a narrow band of intrinsic merits that allows for persistence of an allele, therefore only a negligible proportion of the alleles present had an intrinsic merit below that band (figures B.22, B.23 and B.24 show the exact proportions for a particular threshold in more detail). The DAA model, however, featured a characteristic tail, containing a substantial proportion of alleles in this area of low intrinsic merit (figures 4.13, B.20 and B.21 show the proportions for two particular thresholds). Moreover this characteristic behaviour of both models was insensitive to the population size (figures B.13, B.14, B.15, B.16, B.17 and B.18).

Examining the distribution of allele intrinsic merits shown in figures B.13, B.14, B.15, B.16, B.17 and B.18 demonstrated that figures 4.13, B.20 and B.21 indeed used suitable candidates for allele intrinsic merit thresholds, although higher threshold values may have been appropriate for higher population sizes. However, for reasons of simplicity no such dependency of the threshold value on population size was defined. Increasing the intrinsic merit threshold further would at some point cut into the stable band of allele intrinsic merits, and the goal of separating the two groups of alleles with different behaviour would not have been achieved. Decreasing the intrinsic merit threshold is possible to some extent, until at some point, when the lower end of the tail is reached, the proportion of alleles below that threshold value would become zero.

Having identified suitable thresholds enables me to examine properties of alleles in the low-intrinsic merit zone of the different overdominance models. This ultimately helps in understanding why these alleles are much more frequent in the DAA model than in the asymmetric overdominance models. The first of the properties to be examined was the persistence time of alleles with low intrinsic merit (section 5.3.1.2).

5.3.1.2 Persistence times of alleles with low intrinsic merit

Figures 4.14, B.25 and B.26 show the mean age of the alleles with low intrinsic merit under the DAA and the SAsOD model for different intrinsic merit thresholds. These figures clearly illustrate that the DAA model both featured a much larger proportion of alleles with low intrinsic merit and a much higher average persistence time of these alleles than the asymmetric overdominance models. Even for the stricter intrinsic merit threshold, alleles with low intrinsic merit lasted on average tens of thousands of generations in the DAA model, but generally vanished after just around a hundred generations in the asymmetric overdominance models (figures B.27, B.28 and B.29).

The average allele age in the DAA model only decreased for very large population sizes. This was because, as outlined before, the intrinsic merit threshold used did not account for population size. As higher population sizes resulted in a higher total number of mutations, more alleles with high intrinsic merit were generated. Therefore the mean allele intrinsic merit, and with it the distance of the stable band of alleles from the alleles with low intrinsic merit, increased with population size, resulting in a stronger selection pressure on the alleles with low intrinsic merit. Such a phenomenon, however, was not noticeable in the asymmetric overdominance models in the same way - alleles with low intrinsic merit in these models were “transient” to a very high degree for all population sizes. The next section (section 5.3.1.3) defines this transience of alleles more precisely.

5.3.1.3 Intermediate and transient alleles

Applying another potential measure of allelic diversity, the total range of allele intrinsic merits at the end of the simulation time span, demonstrates the need to distinguish between two different types of alleles with low intrinsic merit. Figures 5.1 and C.1 show the absolute intrinsic merit range, i.e. the range of allele intrinsic merits that incorporated all alleles that existed at the end of the 10 million generations.

Even though figures 5.1 and C.1 illustrate that the allele intrinsic merit range is higher for the DAA model than for the traditional asymmetric overdominance models, the difference was not as high as expected, given the large difference in allele intrinsic merit standard deviation (figures 4.12 and B.19).

This absolute allele intrinsic merit range is however an inappropriate measure particularly for large population sizes, since many of the alleles in the final set had an extremely short life span and only existed at very low frequencies. In the extreme case of very high mutation rates, one would just be measuring a range of allele intrinsic merits close to that used to generate newly mutated alleles, and would therefore be unable to spot differences between the models. Figure 5.2 illustrates this dilemma. It shows a number of alleles that have

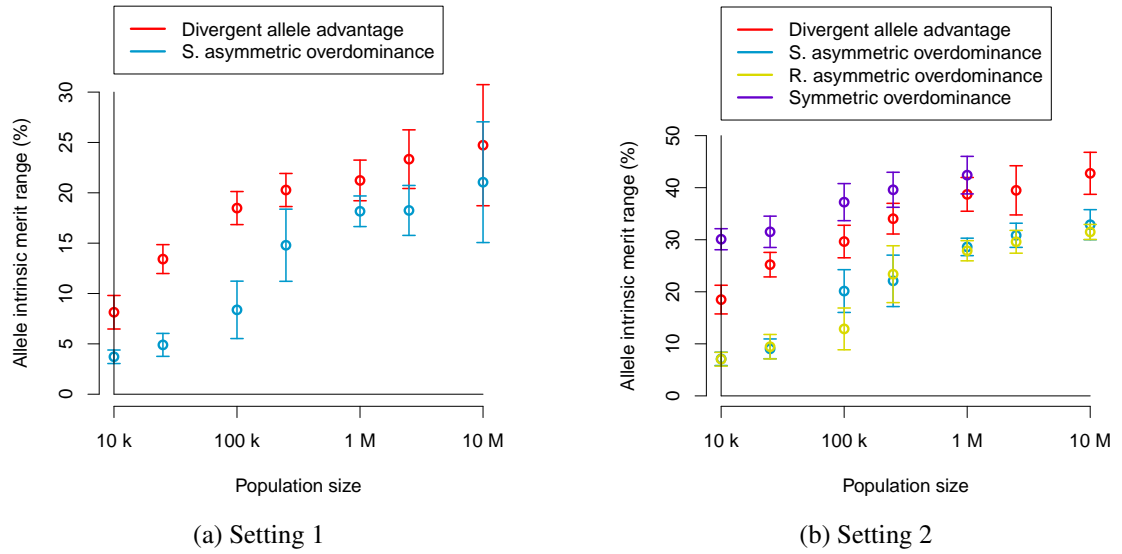


Figure 5.1: Allele intrinsic merit range vs. population size (taking all alleles at the end of the simulation period into account) for settings 1 (top) and 2 (bottom).

just recently emerged, but some of them – particularly in the SAsOD model – were very infrequent and short-lived.

Nevertheless, relatively recent alleles can have a major impact on the gene pool of a population, as is the case for the DAA model. This leads to the conclusion that some kind of threshold is needed to usefully exclude alleles that are either too rare or too recent from the range measure that is applied. Apart from giving a better idea of the effective range of allele intrinsic merits, this threshold measure can usefully be applied to gain more insights into the behaviour of the models by classifying alleles with low intrinsic merit as either “transient” or “intermediate”.

I define alleles with low intrinsic merit (i.e. with an intrinsic merit value below the intrinsic merit threshold discussed in section 5.3.1.1) as “intermediate” alleles if they are either frequent enough or have not emerged too recently. Therefore, “transient alleles” are defined as alleles with very low frequency, or alleles that emerged very recently. This again requires the definition of a threshold, which now refers to the emergence time or frequency of an allele. I subsequently call the threshold based on the emergence time of an allele as the “age threshold”, which excludes alleles that emerged in generations more recent than the threshold value (e.g. an age threshold of 10 means that alleles that emerged in the last generation or the 9 generations before are excluded). The threshold based on the frequency of an allele, the “frequency threshold”, excludes alleles that are less frequent (have a lower number of representatives in the gene pool) than the threshold value.

Figures 5.3 and 5.4 show numbers of alleles and share of the gene pool of intermediate and transient alleles for the DAA model and the asymmetric overdominance models for settings

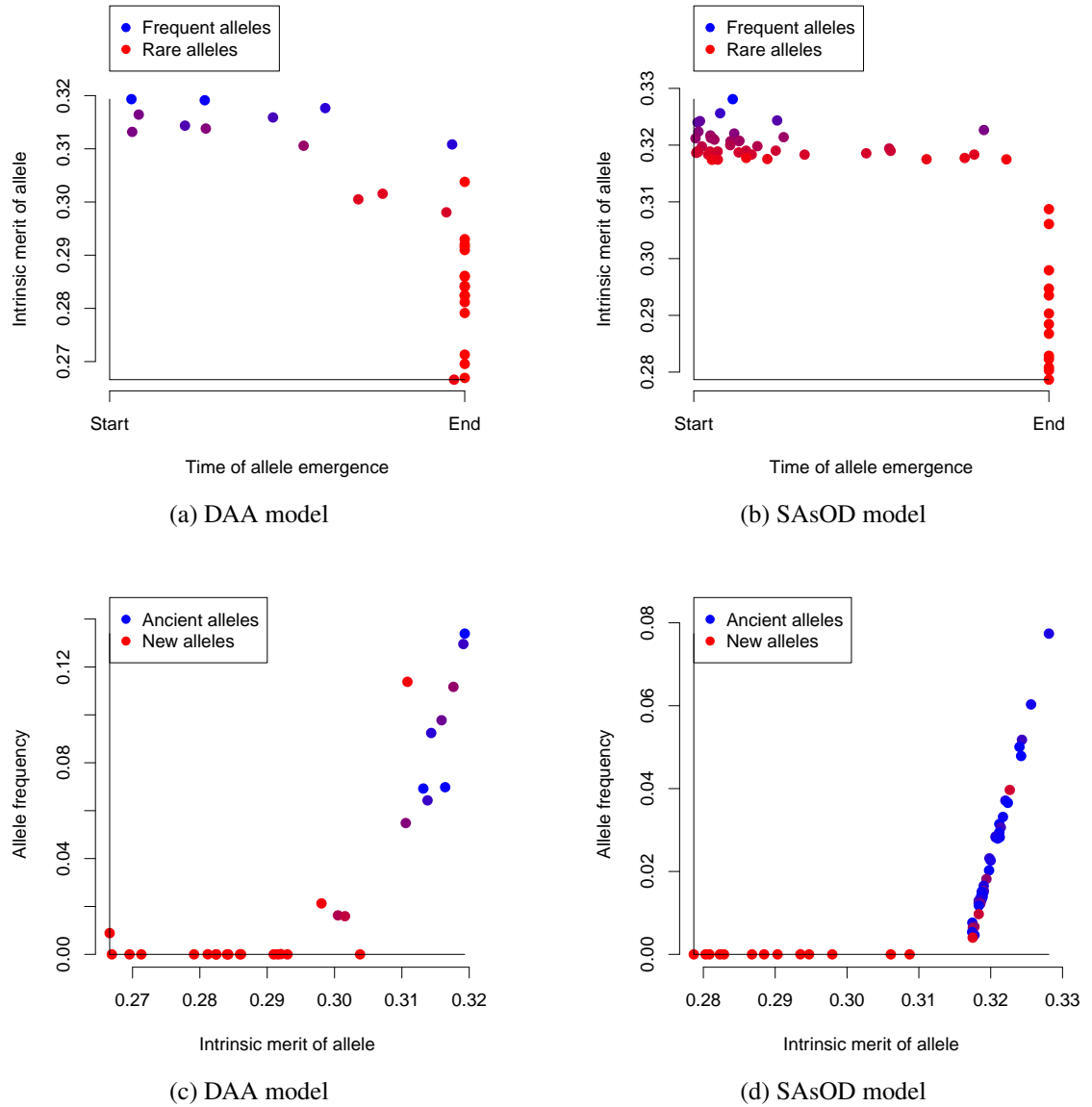


Figure 5.2: Allele intrinsic merit against time of allele emergence (top) and allele frequency against allele intrinsic merit (bottom) of single simulation runs for the DAA model (left) and the SAsOD model (right). The population size was 1 million individuals.

1 and 2 and an intrinsic merit threshold of $w_t = 0.3$ for setting 1 and $w_t = 0.35$ for setting 2, and an age threshold of 250 generations.

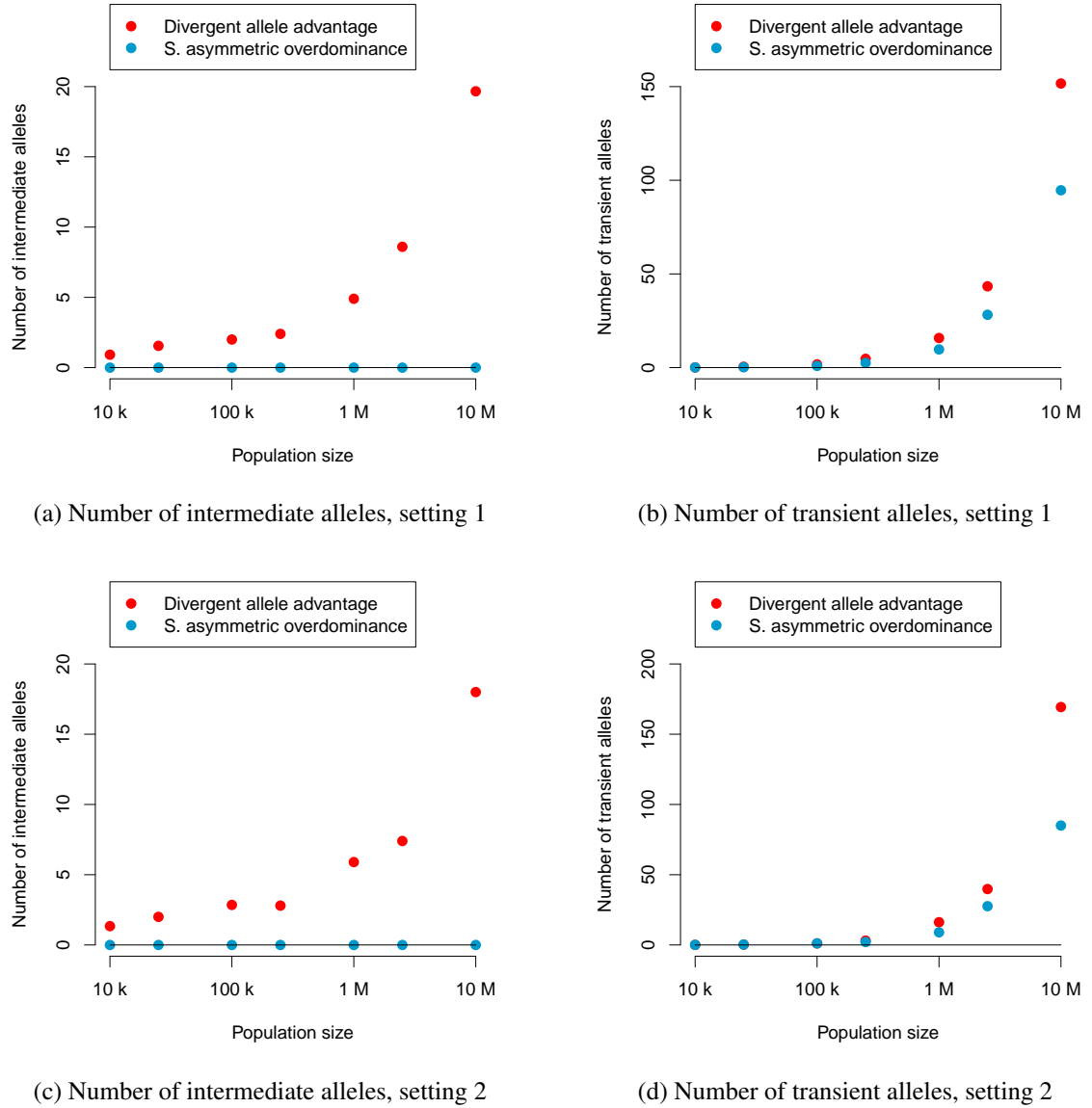


Figure 5.3: Number of intermediate and transient alleles in the DAA and SAsOD models for settings 1 (top) and 2 (bottom). The threshold for allele intrinsic merit was 0.3 for setting 1 and 0.35 for setting 2, and alleles that emerged in the last 250 generations were classified as transient.

They confirm what could be expected given the results illustrated in figure 4.11 and 5.2: There were hardly any intermediate alleles in the tail of the intrinsic merit distributions for the asymmetric overdominance models (in most cases not a single intermediate allele could be identified). In contrast, the alleles with low intrinsic merit made up a comparatively high proportion of the gene pool at the end of the simulation in the DAA model, and intermediate alleles made up the vast majority of this share, although they did not occur in large numbers (only relatively few intermediate alleles were identified for settings 1 and 2 and lower

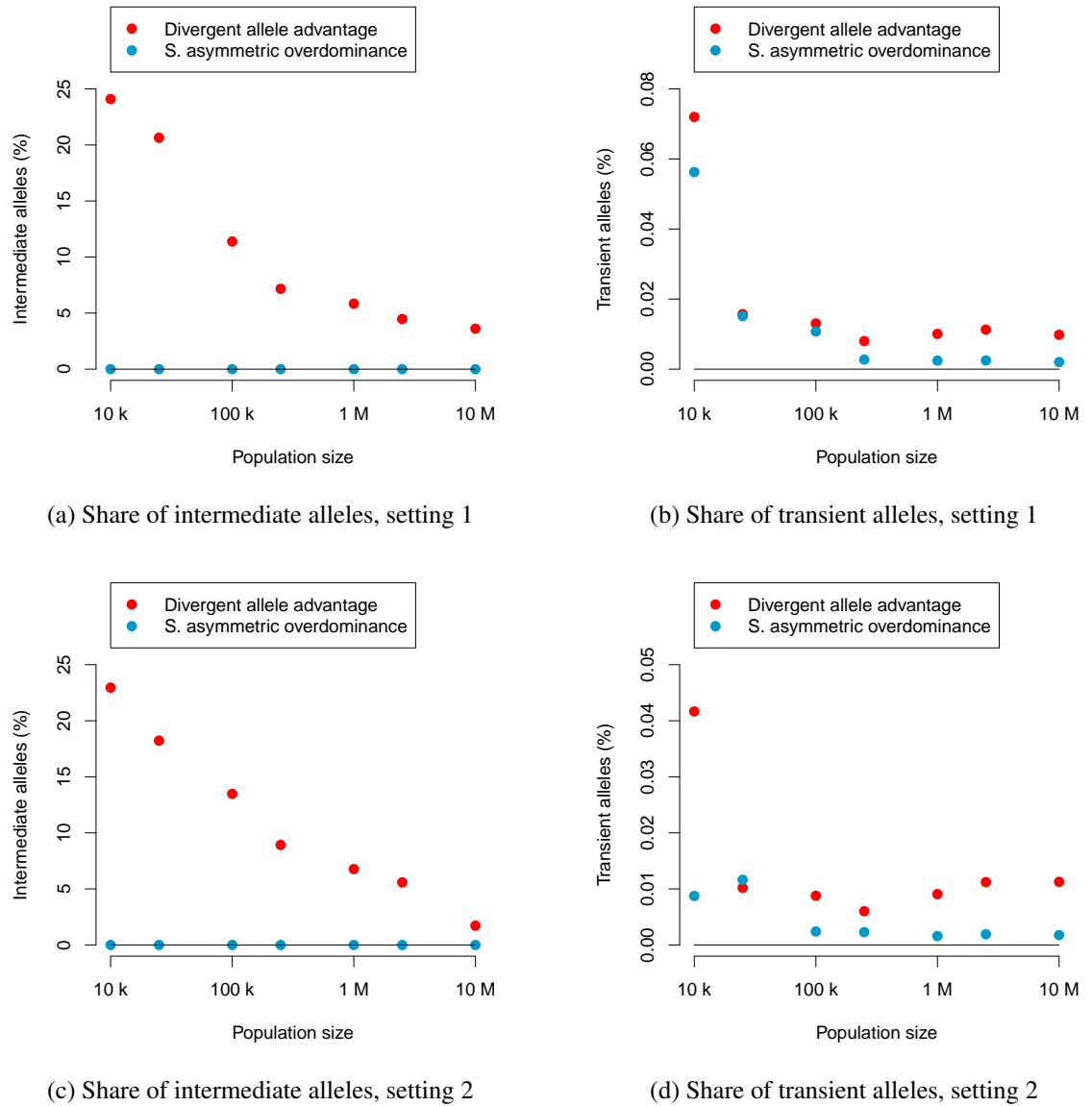


Figure 5.4: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 1 (top) and 2 (bottom). The threshold for allele intrinsic merit was 0.3 for setting 1 and 0.35 for setting 2, and alleles that emerged in the last 250 generations were classified as transient.

population sizes, see figure 5.3).

A more detailed picture regarding the share of the gene pool of intermediate and transient alleles, involving all settings and two different thresholds for allele intrinsic merit, can be seen in the appendix: figures C.9 and C.10 show the share of alleles for the default intrinsic merit threshold, while figures C.11, C.12 and C.13 demonstrate that the same qualitative results were obtained when a stricter threshold for allele intrinsic merit was applied.

The share of the gene pool of intermediate alleles was also not sensitive to the type of threshold used, i.e. the age threshold or the frequency threshold, nor was it sensitive to the exact threshold value. This is demonstrated in figures 5.5 and 5.6, which compare the gene pool share of intermediate alleles for setting 1. Alleles with low intrinsic merit were selected by applying the stricter threshold value in figure 5.5, while figure 5.6 applied the less restrictive threshold value for allele intrinsic merit. To classify alleles as intermediate or transient, both types of thresholds (age and frequency) were used with two different threshold values, respectively.

This insensitivity to the threshold types and values used for classifying alleles with low intrinsic merit as either intermediate or transient simplified the analysis, which was subsequently performed for the age threshold with a threshold value of 250 generations only. Whilst the exact values of the age and frequency thresholds are to some extent arbitrary, the concept of partitioning alleles with low intrinsic merit into two distinct groups is nevertheless useful in highlighting important differences between the DAA model and the asymmetric overdominance models (see in particular section 5.3.1.4).

While the AsOD models did not possess intermediate alleles, and therefore the intrinsic merit threshold was sufficient, applying the age or frequency threshold without any intrinsic merit threshold generally led to the same separation of transient alleles from the band of stable alleles. This indicates that in the traditional asymmetric overdominance models there is a “natural” partitioning of alleles into stable alleles, which are frequent, old and have a comparatively high intrinsic merit, and transient alleles, which are infrequent, have only recently emerged and have a comparatively low intrinsic merit.

When the same thresholds that partitioned the alleles in the asymmetric overdominance models so naturally into two distinct groups were applied to the DAA model, a third group, the intermediate alleles, appeared. Properties of the intermediate and transient alleles in the DAA model were compared to find out what allowed the intermediate alleles to become temporarily so frequent.

First, differences in allele intrinsic merit were explored. Figures 5.7, C.14 and C.15 (less restrictive intrinsic merit threshold) as well as figures C.16, C.17 and C.18 (stricter intrinsic merit threshold) compare the average intrinsic merit of intermediate alleles to transient alleles. These figures clearly demonstrate that it is not intrinsic merit that caused intermediate

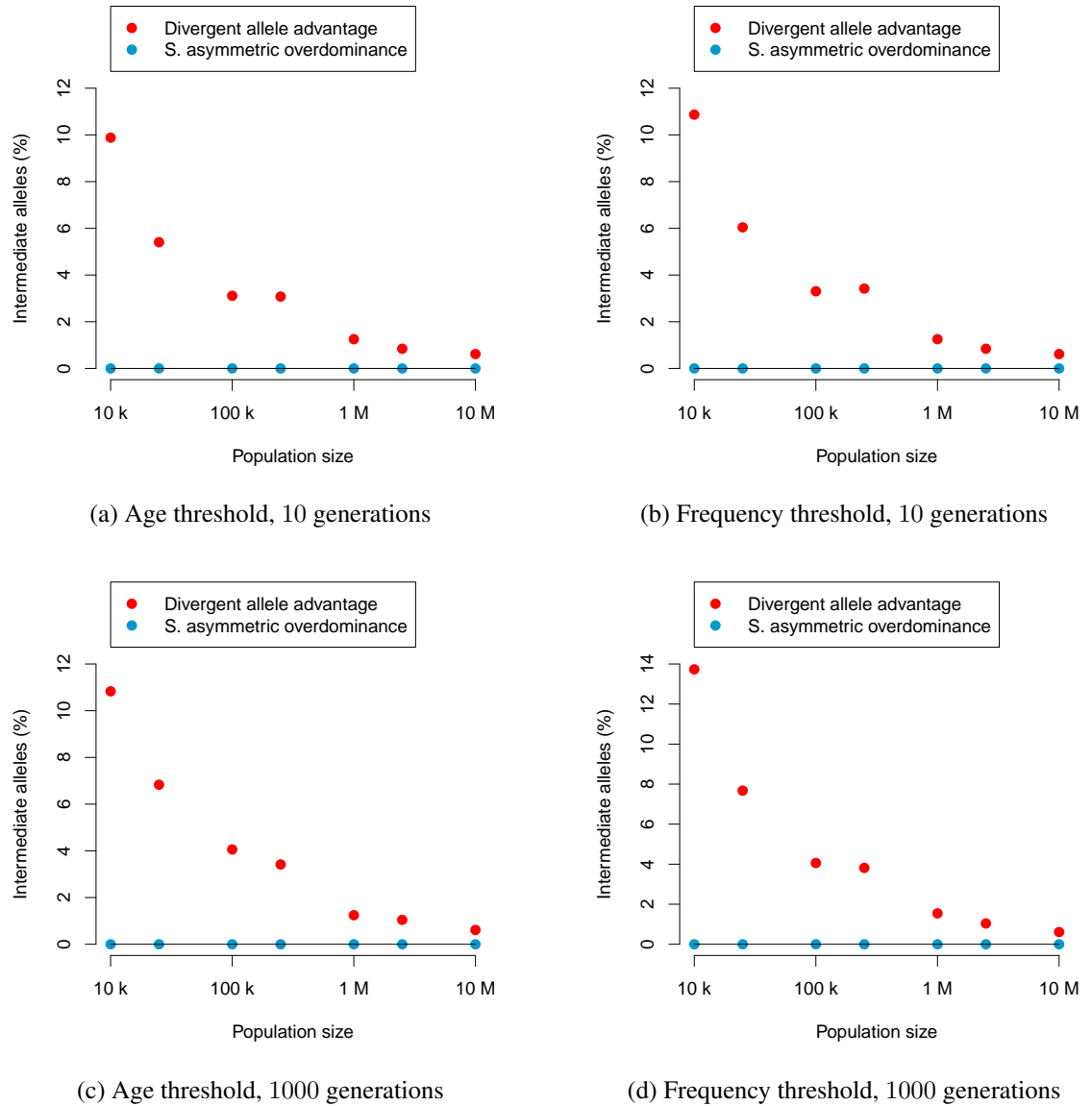


Figure 5.5: Share of the gene pool of intermediate alleles in the DAA and SAsOD models for setting 1 and different age and frequency thresholds. The threshold for allele intrinsic merit was 0.275, which corresponds to the lower value used, i.e. the stricter threshold.

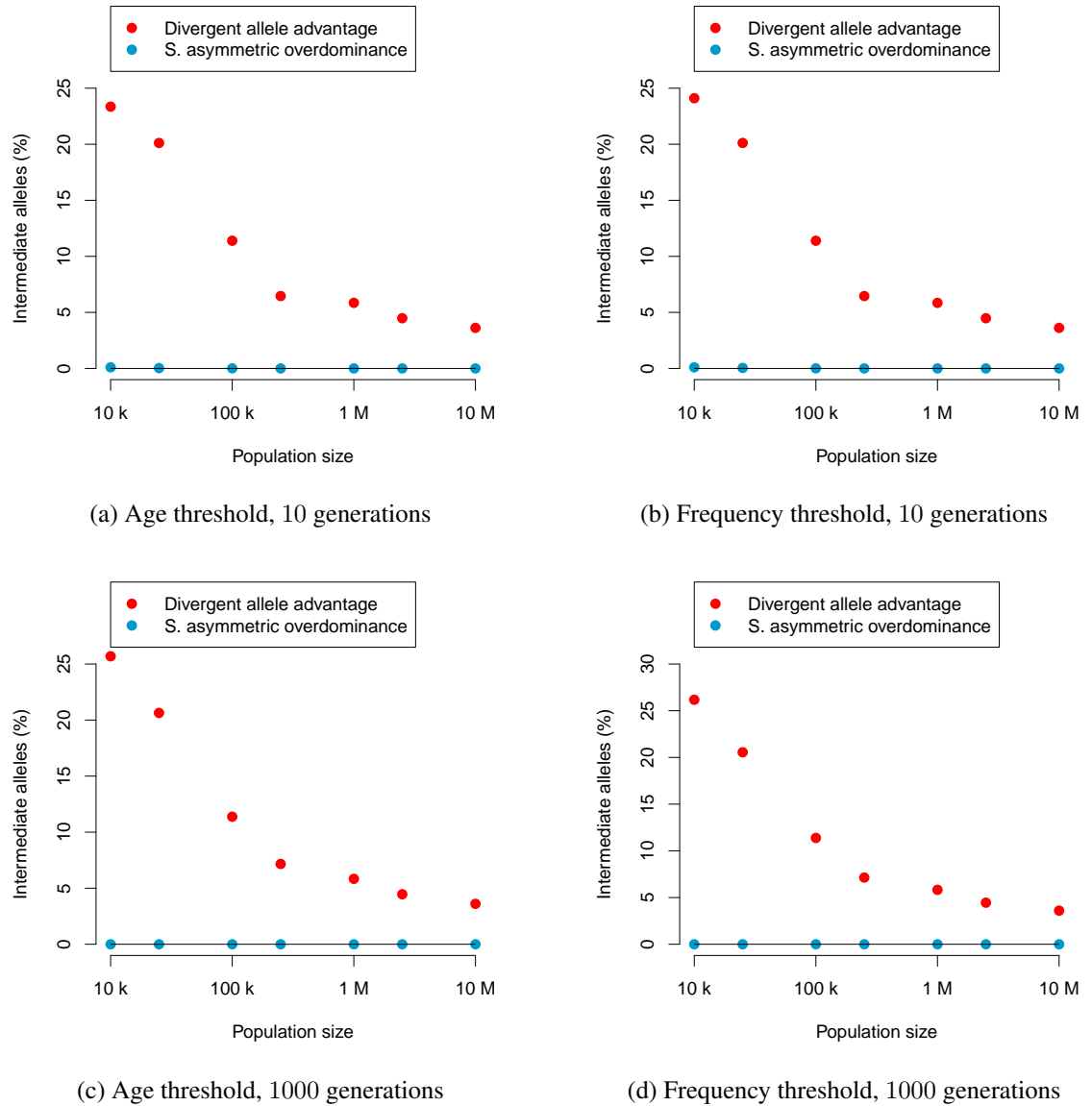


Figure 5.6: Share of the gene pool of intermediate alleles in the DAA and SAsOD models for setting 1 and different age and frequency thresholds. The threshold for allele intrinsic merit was 0.3, which corresponds to the higher value used, i.e. the less restrictive threshold.

alleles to become frequent in the DAA model, as there was no significant intrinsic merit difference between intermediate and transient alleles. Hence, the hypothesis that intermediate alleles are more frequent because they have an intrinsic merit that is on average above that of other alleles with low intrinsic merit can be dismissed.

What did distinguish intermediate alleles from other alleles with low intrinsic merit was the special feature of the DAA model, the importance of the sequence difference between two alleles on the same locus. If the encoded amino acid sequence of an allele was different from the rest of the gene pool, then this allele had a competitive advantage over other alleles with similar intrinsic merit. This selective advantage was responsible for maintaining these strongly different alleles at high frequencies.

I measured how different (in % of variable sites) a particular allele is when compared to any other allele in the gene pool, and used this information to calculate the average sequence difference of this allele to the rest of the gene pool by taking into account all the possible heterozygous genotypes this allele can form as well as their respective frequencies. Due to the model dynamics of the DAA model, which reward alleles for being different, the average sequence difference was much higher in the DAA model than in traditional overdominance models, and it increased with population size. What is more, intermediate alleles had a substantially higher average sequence difference to the rest of the gene pool compared to transient alleles, as figure 5.7 demonstrates. This means that the intermediate alleles distinguish themselves from other alleles with low intrinsic merit by being on average more different to the rest of the gene pool, which helps them persist for longer, and at higher frequencies.

Other traditional asymmetric overdominance models, however, lack the feature of a fitness reward for heterozygotes carrying divergent alleles, and thus do not show the phenomenon of intermediate alleles. A higher intrinsic merit is the only chance for an allele to persist under these models, and this restriction limits the intrinsic merit diversity of the gene pool considerably.

There is of course more than one way to draw a threshold between intermediate and transient alleles. Instead of considering either the emergence time or the frequency of an allele, it would also be possible to combine these two measures, and define alleles that are both recently emerged and very infrequent as transient alleles. Furthermore, frequency or age thresholds could be made dependent on population size to allow for stricter thresholds for smaller population sizes. I will however subsequently only apply simple allele age thresholds, which is sufficient to demonstrate the effectiveness of the concept.

In the next three sections (5.3.1.4, 5.3.1.5 and 5.3.1.6), the definition of transient alleles is expanded to contain all alleles that emerged recently (i.e. that have an age below the threshold value) independent of their intrinsic merit. This is sufficient for clearing the haze of recent mutations that obscured measures as the range of allele intrinsic merits (see figures

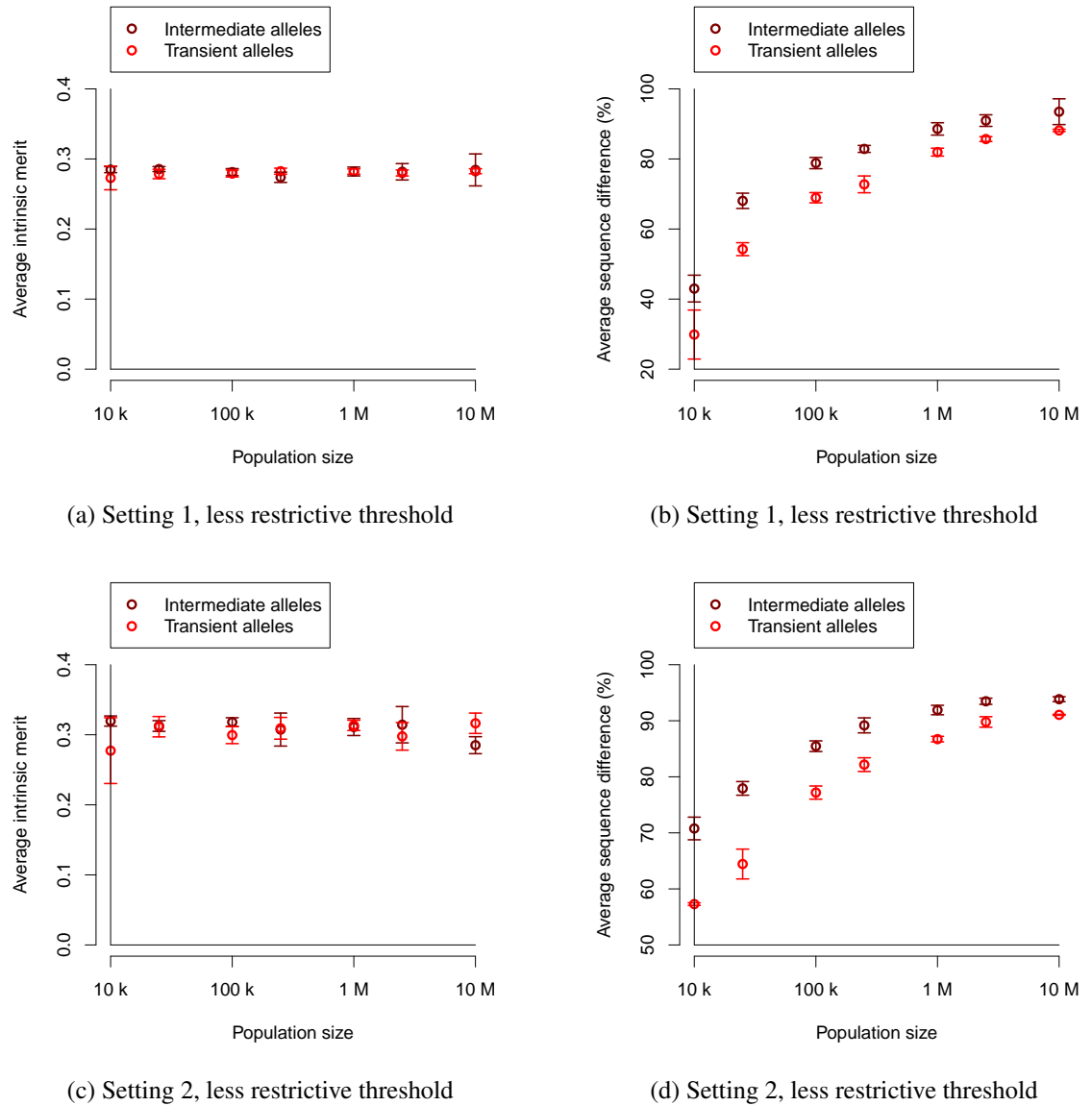


Figure 5.7: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 1 (top) and 2 (bottom). The threshold for allele intrinsic merit was 0.3 for setting 1 and 0.35 for setting 2, and alleles emerged in the last 250 generations were classified as transient.

5.1 and C.1) without arbitrarily modifying the measure of interest. As the way the intrinsic merit is generated for newly mutated alleles (section 4.3.4) does not result in a high intrinsic merit for recent mutations in the vast majority of cases, the difference between the two definitions is minor.

5.3.1.4 Range of allele intrinsic merits

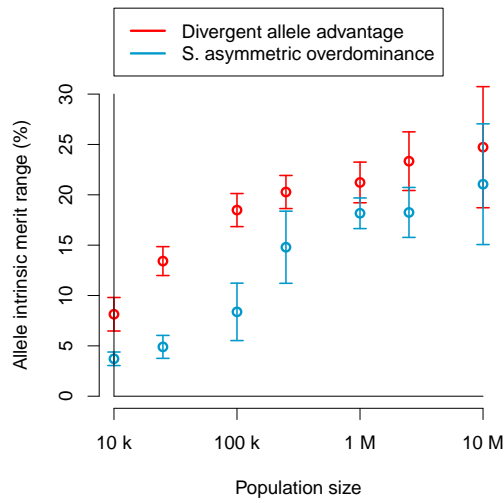
Figure 5.8 shows the range of allele intrinsic merits when recently mutated alleles were filtered out by applying different age thresholds.

This figure demonstrates that excluding at a minimum the alleles in the gene pool that emerged in the last 250 generations led to a fairly stable range of allele intrinsic merits for the asymmetric overdominance models for population sizes of up to 10 million, indicating not only that this threshold is appropriate for excluding transient alleles for the entire range of population sizes investigated, but also that the increase in allele intrinsic merit range for larger population sizes without exclusion of transient alleles is only a consequence of the higher total number of mutations. In the DAA model, however, the range of allele intrinsic merits was only marginally affected by the exclusion of transient alleles, and even an age threshold of 1000 generations (figure 5.8f) did not substantially alter the result when compared with the range of allele intrinsic merits of all alleles (figure 5.8a).

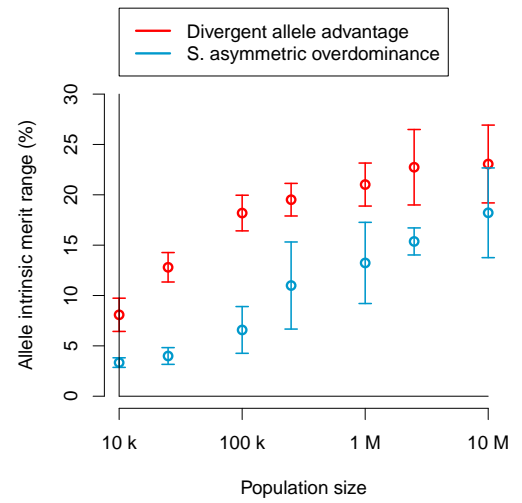
As a consequence, figures 5.8, C.2, C.3, C.4, C.5 and C.6 show that the gap in allele intrinsic merit range between the DAA model (and the SOD model, see figure 5.9) on one hand and the asymmetric overdominance models on the other opened up quite clearly towards large population sizes, starting with a ratio of approximately 3 : 1 for 10000 individuals and increasing to a ratio of approximately 7 : 1 for a population size of 10 million in setting 1, and comparable results for the other settings.

5.3.1.5 Weighted standard deviation in allele intrinsic merit

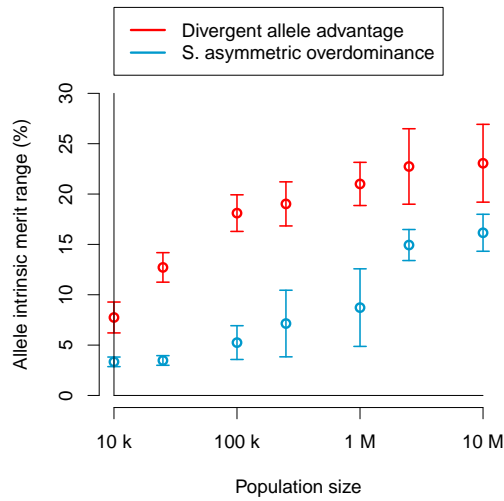
The weighted allele intrinsic merit standard deviation that was used as a measure of diversity for the gene pool of alleles was insensitive to the threshold used (figures 5.10 and C.7), as the transient alleles did not have an appreciable impact on this measure, given their rareness. Consequently, these figures were indistinguishable from figures 4.12 and B.19 where all alleles (including the transient ones) were accounted for. Removing infrequent or recent alleles had little effect on the conclusions derived from the weighted allele intrinsic merit standard deviation, which is by definition essentially independent of infrequent alleles.



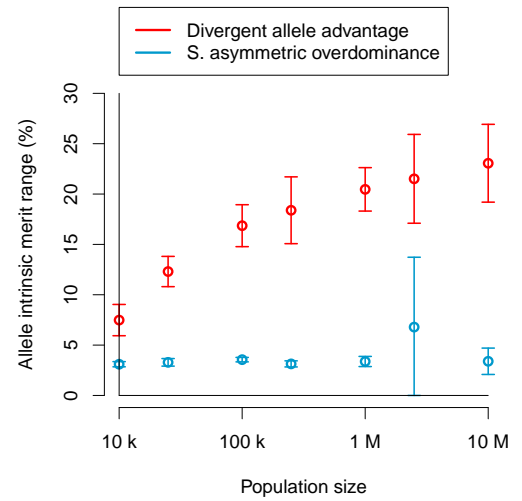
(a) No age threshold (all alleles)



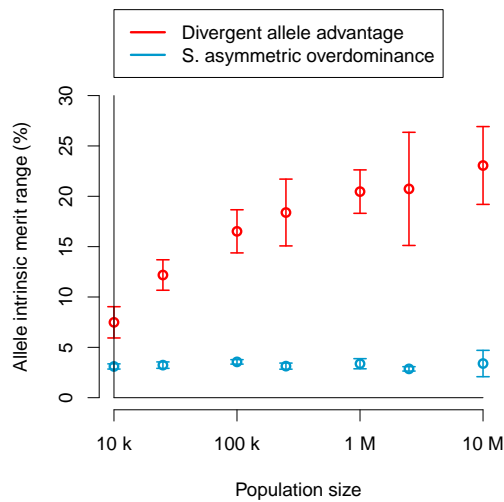
(b) Age threshold 10 generations



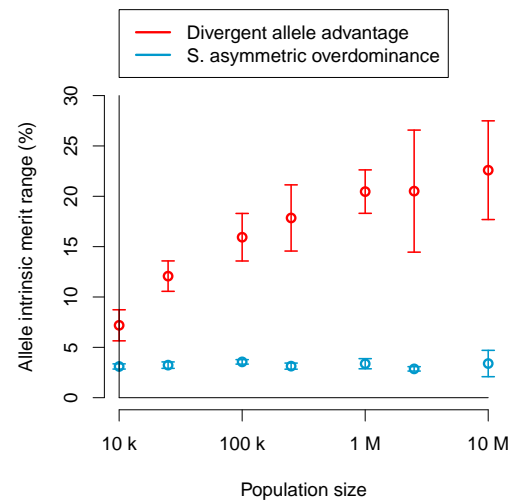
(c) Age threshold 25 generations



(d) Age threshold 100 generations



(e) Age threshold 250 generations



(f) Age threshold 1000 generations

Figure 5.8: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 1.

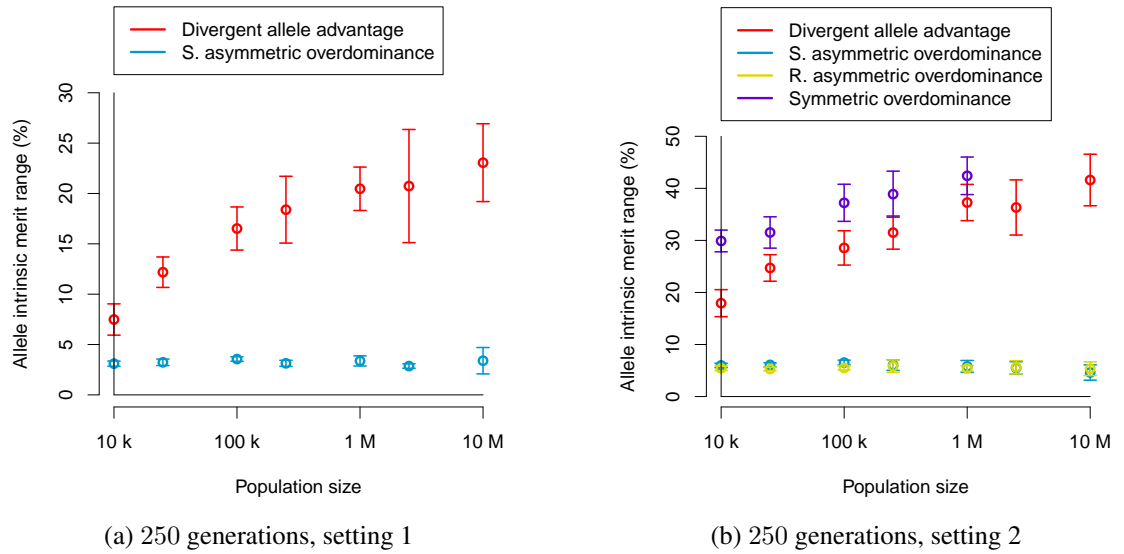


Figure 5.9: Range of allele intrinsic merits of the final set of alleles when using an age threshold of 250 generations to remove transient alleles. Results correspond to settings 1 (left) and 2 (right).

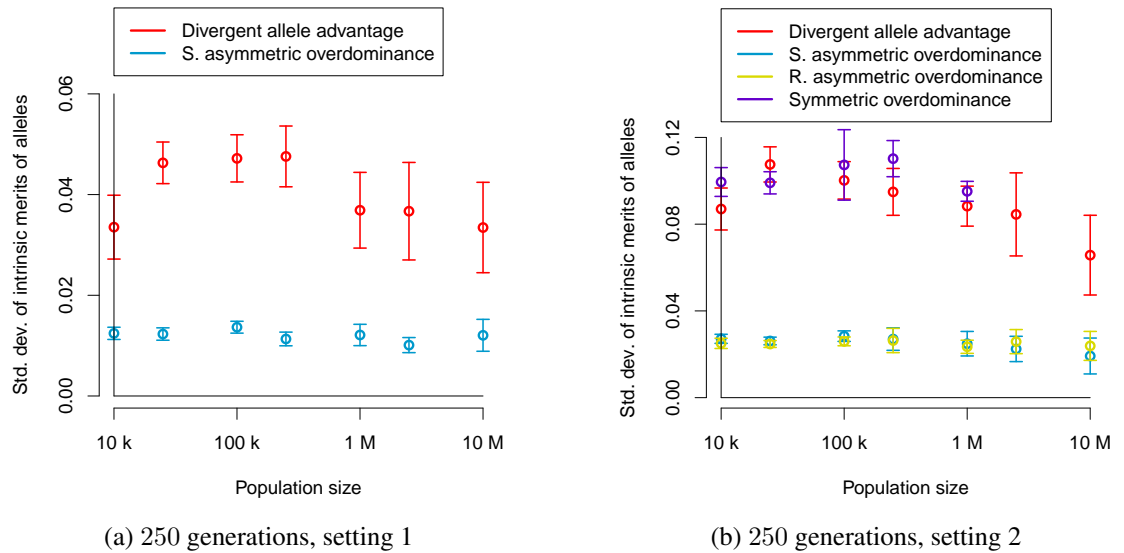


Figure 5.10: Standard deviation in allele intrinsic merit of the final set of alleles when using an age threshold of 250 generations to remove transient alleles. Results correspond to settings 1 (left) and 2 (right).

5.3.1.6 Number of stable and intermediate alleles

Applying a threshold that excludes transient alleles did however have an effect on allele numbers at the end of the 10 million generations of evolution, as demonstrated by figures 5.11 and C.8 when compared to their counterparts without exclusion of transient alleles (figures 4.7 and B.5).

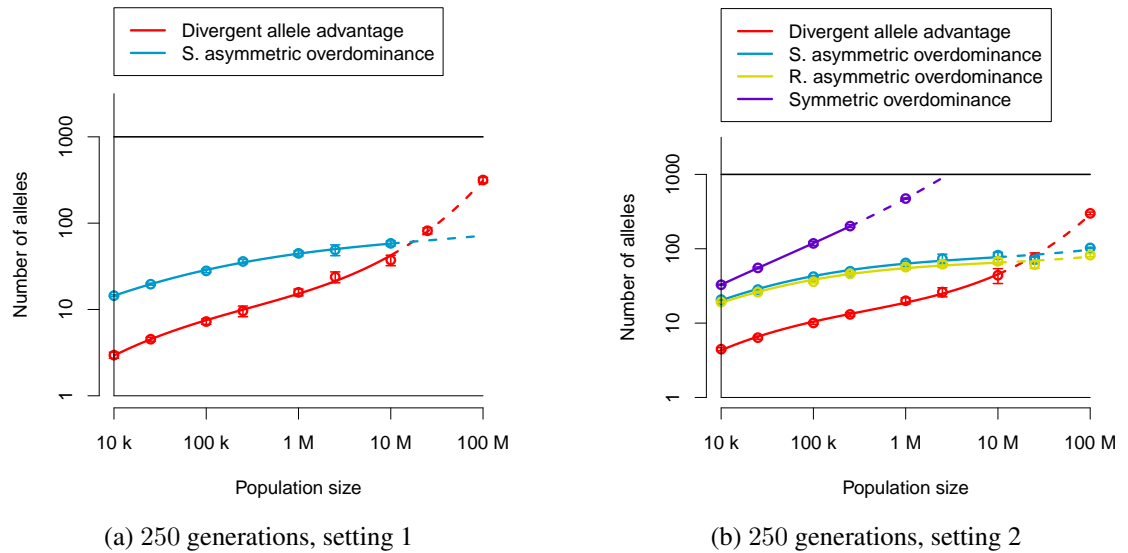


Figure 5.11: Number of alleles at the end of the simulation when using an age threshold of 250 generations to remove transient alleles. Results correspond to settings 1 (left) and 2 (right).

Therefore, despite an increase in allele numbers for larger population sizes in the case of the DAA model, the range of intrinsic merits of non-transient alleles increased, whereas for the asymmetric overdominance models the number of non-transient alleles leveled off for large population sizes, and ranges of allele intrinsic merits stayed constant for all population sizes. This confirms the findings of De Boer et al. (2004) that the asymmetric overdominance models are only able to show “a high degree of polymorphism if, and only if, the fitness contributions of MHC alleles are very similar”, but also demonstrates that this effect is not a general feature of heterozygote advantage, but rather of the specific model applied: traditional asymmetric overdominance models exhibit this behaviour, but the DAA model does not.

5.3.2 Population fitness and alleles with low intrinsic merit

While alleles with low intrinsic merit produce poor homozygotes, their presence in the gene pool of a population can still improve the overall population fitness, provided these alleles have an above-average sequence difference to the most frequent alleles, and thus form fit

heterozygotes. For that reason, a marked decrease in population fitness occurred in the DAA model when alleles with low intrinsic merit were removed, while traditional overdominance models showed a substantially lesser decrease. This is demonstrated in figure 5.12, which shows scenarios where the 25% or 50% of the allele representatives with the lowest intrinsic merits have been removed for setting 1 and 2, respectively.

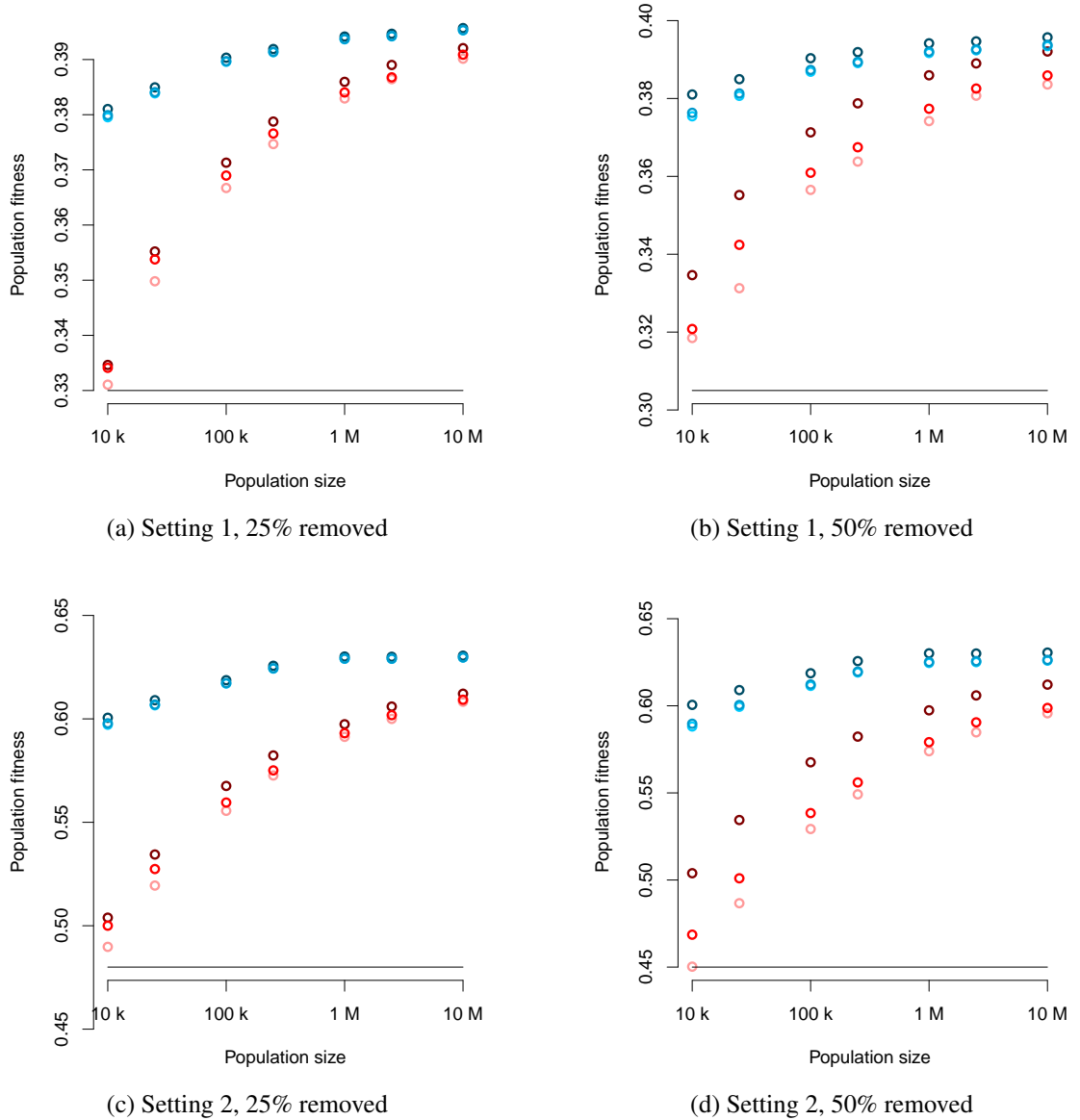


Figure 5.12: Population fitness for the DAA model (shades of red) and the SAsOD model (shades of blue). The dark circles correspond to the population fitness of a population of given size after 10 million generations, while the light circles correspond to the population fitness of the same population, but with a certain share of the allele representatives with the lowest intrinsic merits (25% or 50%) removed. The medium circles correspond to the same population as the light circles, but with optimal frequencies of the remaining alleles.

5.3.3 Population Fitness during and after a severe population bottleneck

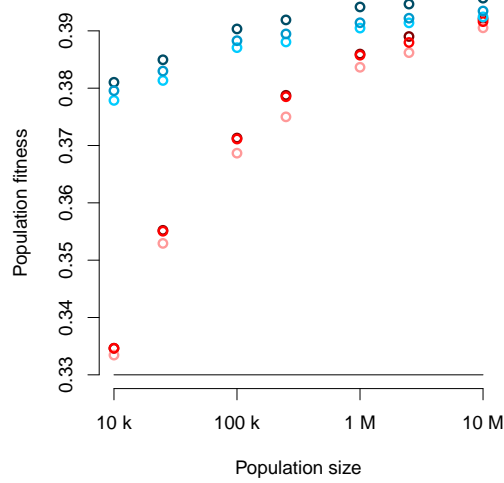
A similar method can also be applied to explore changes in population fitness during population bottlenecks. Figures 5.13 and 5.14 show the changes in population fitness during and after bottlenecks of different severities for settings 1 and 2. The bottleneck was modelled by considerably reducing the population size during a short time span (the bottleneck phase, see table 5.1 for the different severities of the bottleneck). If the population size was reduced to m_b individuals during the bottleneck, then $2m_b$ alleles were randomly drawn from the gene pool of the final population (i.e. the gene pool at the end of the simulation span of 10 million generations). After that, the population size was assumed to recover quickly (in the simulation this happened immediately) to its initial level. The population fitness was recorded throughout the bottleneck in figure 5.14.

scenario		1	2	3	4
number of individuals	DAA	100	10	5	3
	AsOD	100	10	5	3
average number of alleles	DAA	9.95	7.25	6.35	4.1
	AsOD	40.45	15.75	8.85	5.45
reduction in allele numbers	DAA	13%	36%	44%	64%
	AsOD	7%	64%	80%	87%

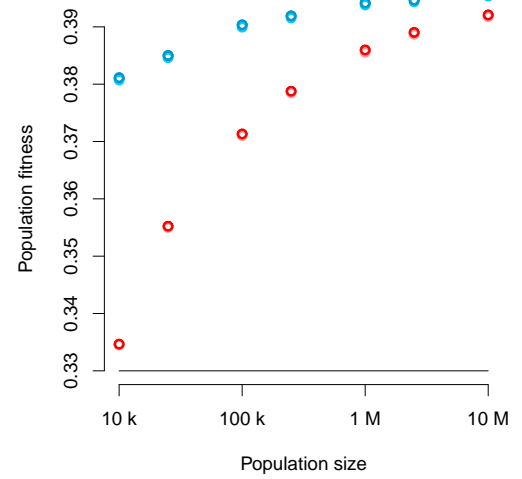
Table 5.1: Population bottlenecks of four different severity levels for the DAA model and the SAsOD model. The initial size of the population is 100000 for all scenarios. Different severity levels correspond to different shades of red (DAA model) and blue (SAsOD model) in figure 5.14. The more severe the bottleneck, the lower the population fitness during the bottleneck, and the longer the duration of recovery.

While the population size during the bottleneck (m_b) is clearly unrealistically low for real-world scenarios, table 5.1 additionally lists the corresponding allele numbers during the bottleneck and the relative contraction of allele numbers (compared to the number of alleles before the bottleneck) for both the DAA model and the SAsOD model, as these numbers can be compared to observational data more usefully.

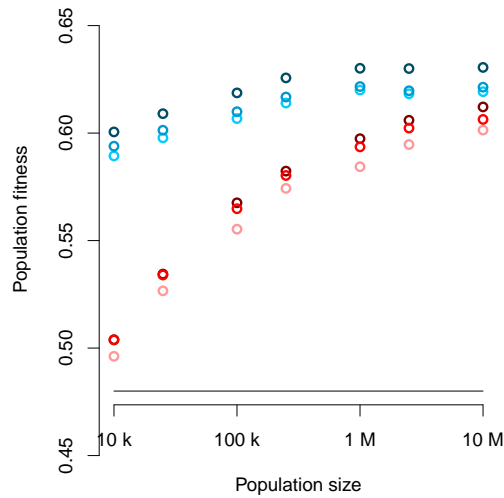
The DAA model exhibited a different response to such a bottleneck compared to traditional models of asymmetric overdominance. While extremely severe bottlenecks led to a slow recovery of population fitness under both models, the preservation of a slightly higher number of alleles helped the population recover a substantial proportion of their fitness much quicker in the DAA model (figure 5.13 shows this for different population sizes). This was because



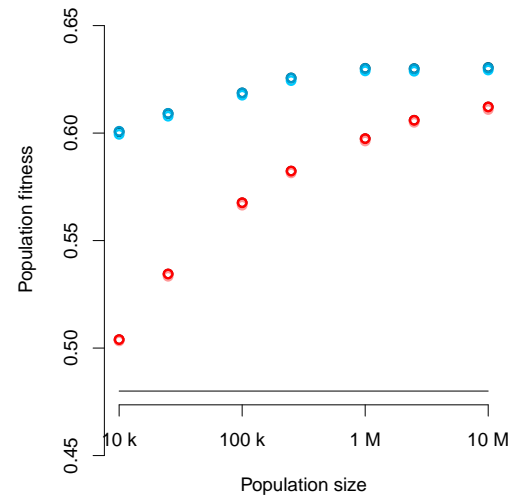
(a) Setting 1, 10 individuals



(b) Setting 1, 100 individuals



(c) Setting 2, 10 individuals



(d) Setting 2, 100 individuals

Figure 5.13: Population fitness for the DAA model (shades of red) and the SAsOD model (shades of blue). The dark circles correspond to the population fitness of a population of given size after 10 million generations, while the light circles correspond to the population fitness in a severe population bottleneck where just a low number of individuals (10 or 100) survived. The medium circles correspond to the same population as the light circles, but with optimal frequencies of the remaining alleles.

the frequencies of the remaining alleles adjusted quickly so that population fitness was optimised for the reduced set of alleles. For this process, new beneficial mutations were not required. However, such a frequency adjustment only led to a marked increase in population fitness in the DAA model when main allelic lineages were preserved.

In the SAsOD model on the other hand, the population fitness increase in that initial adjustment phase was lower: under the DAA model, population fitness reached its original level after around 40% of the time that it took the SAsOD model to restore population fitness when the number of individuals was reduced to 10 during the bottleneck, and the setting was 1 (figure 5.14, second darkest shade). The subsequent rate of population fitness increase was, on the contrary, higher in the SAsOD model compared to the DAA model. This was because the latter required not only new mutations which give rise to alleles with high intrinsic merit, but also the buildup of sequence divergence, which is – as figures 4.3, B.2, B.3 and 5.14 demonstrate – a relatively slow process.

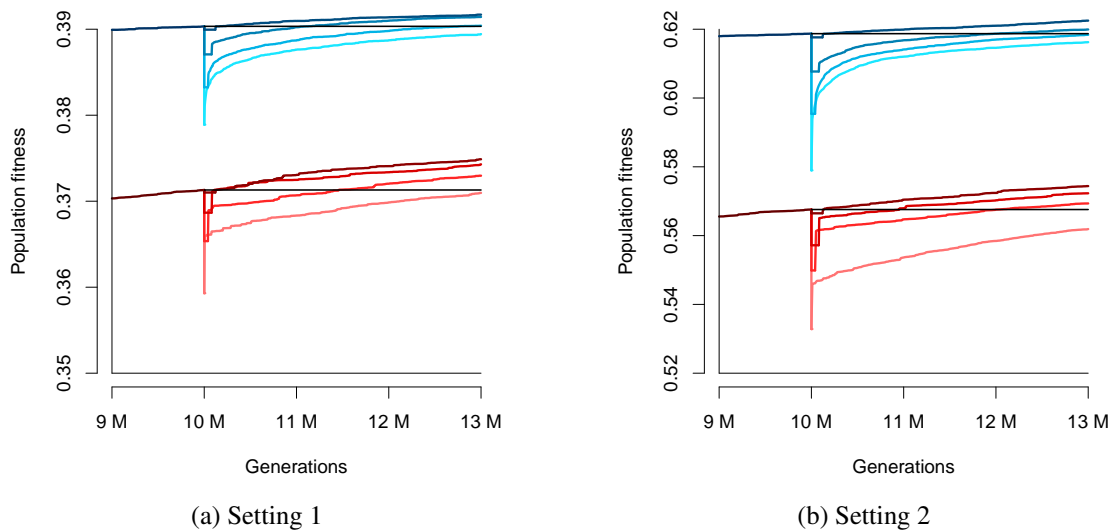


Figure 5.14: Population fitness for the DAA model (shades of red) and the SAsOD model (shades of blue) during a population bottleneck. The number of individuals during the bottleneck varied from 100 (darkest shade) to 10, 5 and 3 (lower numbers of individuals are represented by lighter shades). The population size outside the bottleneck was 100000.

5.4 Discussion and Conclusions

In this chapter I discussed a number of implications that emerge under my novel model based on the divergent allele advantage hypothesis, which was introduced in chapter 4.

The analysis revealed a substantial negative side effect of MHC diversity, which renders several groups of individuals relatively susceptible to disease, as DAA creates more individuals

with more susceptible MHC genotypes than traditional models of heterozygote advantage. Homozygotes are more susceptible under all overdominance models but the proportion of homozygous individuals is higher in the DAA model than in the asymmetric overdominance models (section 4.4.2.4 and figure 4.15). Heterozygotes with similar alleles are also relatively susceptible to disease in the DAA model but not under traditional asymmetric overdominance models. In addition, alleles of low intrinsic merit, which occur at substantial frequencies in the DAA model (sections 4.4.2.3 and 5.3.1.1 and figures 4.11 and 4.13), can produce susceptible heterozygotes and very susceptible homozygotes. These alleles will be maintained in a population under DAA if they form favourable interactions with other alleles to produce relatively fit heterozygotes. Such compensation is not possible in traditional overdominance models (De Boer et al., 2004), where there is a strong correlation between intrinsic merit and fitness averaged over all genotypes. All these relatively susceptible individuals create a genetic burden which is greater than previously recognised. This genetic burden helps explain the persistence of substantial variation within natural and managed populations in susceptibility to disease.

Modelling the evolution of MHC diversity suggests that a necessary by-product of increased population health is the increased susceptibility of particular individuals; however, empirical verification of this prediction is still outstanding. The ability to identify MHC genotypes that increase susceptibility to infectious and parasitic diseases might simplify personalised medicine and could make it easier to focus resources on individuals at increased risk of infection. Increased risk could arise from higher levels of exposure as well as compromised immune responses. If experimental evidence supports the usefulness of the proposed model based on the DAA hypothesis, then selecting animals with a set of highly divergent alleles and optimising the allele frequencies to reduce the number of unfit individuals would improve disease resistance in livestock, or in managed populations of wild species.

It was demonstrated that a relatively large proportion of alleles (counting the number of allele representatives) has a low intrinsic merit in the DAA model. This means that, despite selection forces being at work, a substantial part of a population may carry at least one allele that is not very fit. Identifying these alleles and potentially excluding them from being passed on to homozygous offspring, or alternatively selecting a set of sufficiently divergent alleles in breeding programmes might support improving disease resistance in livestock, or in wild species of conservation concern, always assuming that experimental evidence for the divergent allele advantage hypothesis accrues.

Finally, I showed that under the DAA model, the time it takes for a population to recover its fitness after a crash in numbers strongly depends on the severity of the bottleneck. If the bottleneck is not extremely severe, then population recovery might be faster than previously anticipated, as the most frequent alleles adjust to the optimal frequencies for the reduced set of alleles much quicker than new mutant alleles can establish themselves, which is the

only way traditional overdominance models can restore population fitness. However, if the bottleneck is particularly severe, then population fitness can only recover slowly. Under DAA, the evolution of a system of highly divergent alleles necessarily takes a long time, and until levels of diversity (in terms of amino acid sequence difference) have finally recovered, population fitness will remain at substantially lower levels, potentially putting a population of conservation concern at risk of increased disease susceptibility. These results might help to address the compelling need for a better explanation of the consequences of MHC diversity for population viability (Radwan et al., 2010).

All these observations together can significantly improve our understanding of the MHC, and open the way for better management of human health and better control of disease in wild and managed animal populations if the model predictions can be empirically verified.

Chapter 6

MHC Diversity in Well-mixed and Structured Populations

6.1 Introduction

Chapters 4 and 5 discussed results obtained from stochastic, agent-based simulation models of populations subject to the evolutionary processes of mutation, selection and genetic drift. These models assumed single, well-mixed populations, therefore one of the main evolutionary processes, gene flow caused by migration, was not accounted for. This chapter relaxes the assumption that populations are well-mixed, and explores a metapopulation of (well-mixed) subpopulations that are connected via migration. This is of particular significance for large populations, which are more likely to have some degree of spatial structure.

Examining the creation and maintenance of MHC diversity in a metapopulation pursues two objectives. First, it aims to investigate whether the main results presented in chapter 4 are still valid in this more realistic setup. If true, then the case for the DAA model as the main driver of diversity at the MHC would be strengthened. To answer this question, an extensive sensitivity analysis is performed, varying the migration rate over up to 11 orders of magnitude (depending on population size), with four different degrees of subdivision (by varying the number of subpopulations) and two spatial arrangements that result in different levels of connectivity of the subpopulations. Second, it explores the question whether subdivision of a population increases or decreases allelic diversity, and how different rates of migration and different spatial arrangements influence the MHC polymorphism.

6.2 Materials and Methods

The same model was used as described in detail in chapter 4, with the addition of structuring the population into a metapopulation of connected subpopulations and the necessary migration mechanisms in the structured population. Total (meta-)population sizes of 10000, 100000, 1 million and 10 million were explored, while the default values for all other parameters discussed in chapter 4 (see table 4.1) were used.

All simulations started with the same initial conditions used in chapter 4, and only a single allele with an intrinsic merit of 0.25 was assumed to exist at the start of the simulations. However, after 5 million generations (i.e. half the total simulation span), the population was split into a certain number of subpopulations. In most cases, the number of subpopulations was 5; only in cases where the sensitivity to the number of subpopulations was explored, were populations consisting of 2, 10 and 25 subpopulations additionally modelled. The number of subpopulations after the split of the initial, well-mixed single population remained constant until the end of the simulation (i.e. for another 5 million generations).

6.2.1 Split of the well-mixed population into subpopulations

To split the well-mixed single population into l subpopulations, a weighted sample was drawn without replacement from the single population, with the weights being the number of representatives of each allele in the gene pool of the single population. As many representatives of an allele were drawn without replacement as existed in the gene pool before the split (i.e. if the population size of the entire population was 10000, then the drawn tuple of alleles consisted of 20000 elements, each of which is a representative of a certain allele).

The elements of the drawn tuple, i.e. the drawn allele representatives, were then assigned to the subpopulations sequentially in a way that the i^{th} subpopulation received allele representatives corresponding to tuple indices $i, i + l, i + 2l, \dots$ ($1 \leq i \leq l$). This procedure resulted in l subpopulations of equal size.

6.2.2 Migration

Migration between the subpopulations was controlled by the migration rate ν , which was defined as the number of allele representatives that were transferred from one subpopulation to another subpopulation in a specified period of time. The minimum number of allele representatives that can be transferred was 2, which corresponds to one individual. A simple setup that assumed equal migration rates between any two subpopulations and in both directions was used, i.e. migration did not alter the size of any subpopulation. Furthermore, I assumed that the migration rates were the same between all subpopulations where migration occurred.

Most of the simulations corresponded to situations with maximum connectivity between subpopulations, such that every subpopulation was connected to every other subpopulation. Such an arrangement was subsequently called a star arrangement of subpopulations. A linear arrangement of subpopulations is another setup that was analysed, as such an arrangement represented a metapopulation with a much lower connectivity of the subpopulations. Linearly arranged subpopulations are not uncommon in nature, for example in riverine systems. In these linear arrangements every subpopulation was connected to just one adjacent subpopulation on either side, except for the two terminal subpopulations, which were only connected to the one adjacent subpopulation. These two scenarios were subsequently called the “star” and the “linear” arrangement.

The migration rate ν was varied over a wide range to represent a large number of real-world scenarios. For a metapopulation of 5 subpopulations, the maximum migration rate was set to 10% of all individuals emigrating from each subpopulation to each other connected subpopulation in every generation. This was subsequently decreased by an order of magnitude until a mutation rate of only 2 individuals (i.e. 4 alleles) per generation was reached. This mutation rate was further decreased so that 2 individuals migrated only every 10, 100, 1000, 10000 and 100000 generations. The final scenario explored was a population consisting of five entirely separate subpopulations with no migration occurring between them.

6.2.3 Diversity profiles of the metapopulations

The diversity profiles ${}^qD^Z$ were calculated using the method proposed by Leinster and Cobbold (2012). Naive diversity was calculated using the identity matrix I as similarity matrix Z , while two different similarity measures were used when calculating similarity-sensitive diversity. The first similarity measure was the difference in intrinsic merit of alleles, calculated according to equation 3.2, while the second similarity measure was the difference between the two encoded amino acid sequences on the considered locus according to equation 4.12.

As this chapter is mainly concerned with metapopulations, it is important to define diversity for the different viewpoints in the context of a population consisting of (more or less connected) subpopulations.

Generally, there are two different ways to account for diversity in a metapopulation. The first possibility is to consider diversity in the entire population by regarding the metapopulation as a single, well-mixed population before calculating the diversity profile. The second possibility is to calculate diversity measures for each subpopulation separately, and then average over all subpopulations. When these two possibilities are applied to calculate diversity profiles ${}^qD^Z$, they result in different values for diversity. For example, if considering diversity

for $q = 0$, which is equivalent to the absolute number of alleles without considering their frequencies (“allele richness”), then the first approach will count the total number of alleles in the entire population, whereas the second approach will count alleles for each subpopulation separately and then calculate the arithmetic mean, resulting in a lower diversity for $q = 0$ if at least one allele does not occur in at least one subpopulation.

A similar situation arises for $q = 2$, which is the viewpoint that is mapped to the inverse of the expected homozygosity of the population in the diversity profile of a well-mixed population. In this case the first approach will return a value for diversity at $q = 2$ that is in general not the inverse of the expected homozygosity of the metapopulation, but rather the inverse of the expected homozygosity of a population that merges all subpopulations into a single population, which is considered well-mixed. The second approach, which calculates the diversity for $q = 2$ for each of the subpopulations separately, and subsequently takes the average, will give a better estimate of the expected homozygosity in the population, given that information about population structure is available. This is because the subpopulations themselves are well-mixed, while the metapopulation as a whole is not. However, as soon as at least a minimum amount of migration was present, both approaches gave similar results.

All diversity profiles were calculated using both approaches, but the results presented in this chapter only refer to the first approach that merges the gene pools of the subpopulations into a single pool. This is because the first approach allows for better comparison between the well-mixed and the structured populations, and structured populations with different numbers of subpopulations, as it does not depend on the size of the subpopulations. Furthermore, the first approach is more generally applicable when drawing comparisons with observed data, as information about structuring of a population is often not available. Therefore, whenever a diversity profile was calculated from observed allele frequencies, the final gene pools of all subpopulations of the metapopulation were combined into a single population, which was then used to calculate the diversity profile.

6.2.4 Overview of the simulations performed

Table 6.1 gives an overview of the parameter settings that were used as input for the sensitivity analysis of migration rate and migration interval. For every spatial arrangement and population size, one simulation was performed for a single, well-mixed population (first line), and one simulation was performed for a metapopulation of isolated subpopulations (last line). All combinations of arrangement / model and total population size were combined with the pairs of values characterising migration in the last two columns, resulting in a total of 138 simulations for this sensitivity analysis (10 different migration interval / migration rate pairs exist for a population size of 10000, 11 for a population size of 100000, 12 for a population size of 1 million and 13 for a population size of 10 million). From the six

num. of subpopulations	total population size	migration interval	migration rate
2	1 million	1	2 individuals
		10	2 individuals
5		100	2 individuals
		1000	2 individuals
10		10000	2 individuals
		100000	2 individuals
25		no migration	no migration

Table 6.2: Schedule of the simulations performed for the sensitivity analysis of the number of subpopulations in the metapopulation. The migration interval is given in generations.

6.3 Results

6.3.1 Number of alleles

As in chapter 4, the first diversity measure explored was the number of alleles that are maintained in a population. Figure 6.1 shows the number of alleles over time for a well-mixed population (shades of red) and for a metapopulation consisting of 5 subpopulations (shades of grey) that emerged from the single population after five million generations. Both scenarios are shown for different population sizes; in case of the metapopulation, each of the 5 subpopulations has a population size of $\frac{1}{5}$ of the original population. When counting allele numbers, the metapopulation of subpopulations is considered as a single population in a sense that the total number of different alleles in the whole population is counted rather than the average number of alleles in a single subpopulation. The same procedure is applied when calculating the range of allele intrinsic merits and the weighted standard deviation in allele intrinsic merit (sections 6.3.2.1 and 6.3.2.2).

Figure 6.1 demonstrates that after the split of the single population into subpopulations, the number of alleles can either increase (figure 6.1a) or decrease (figure 6.1b), depending on the level of migration. If there is no migration between the subpopulations, as in figure 6.1a, allele numbers increase, since the set of alleles diverges between the subpopulations as time passes. However, if there is a small amount of migration, as in figure 6.1b, then allele numbers generally decrease.

Figure 6.2 gives a more detailed summary of the interplay between migration rate, population size and number of alleles maintained after 10 million generations. Migration rates are given as the number of generations between migration events and the proportion of individ-

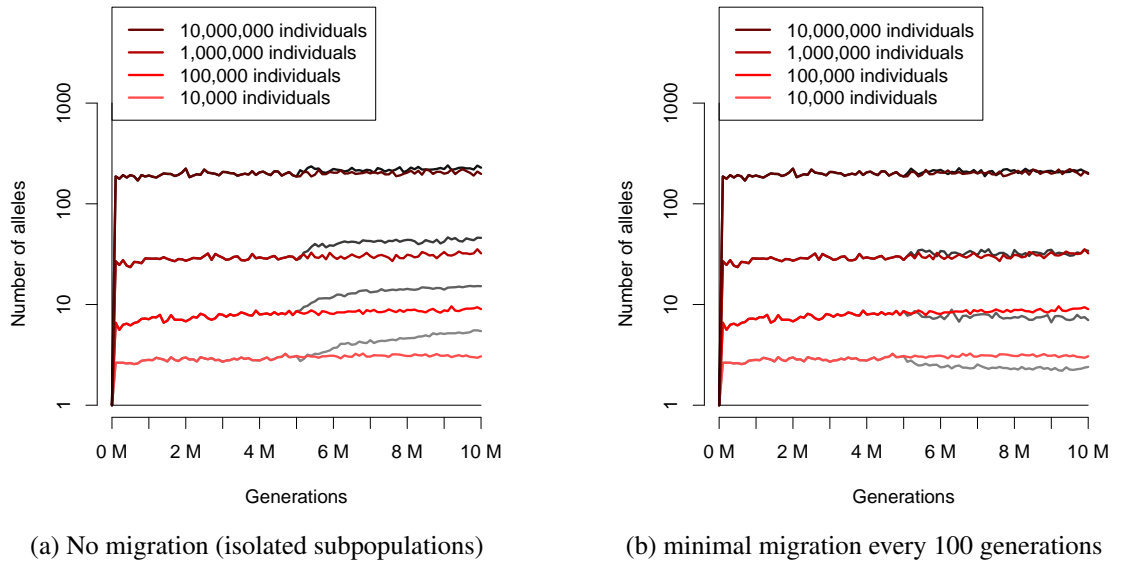


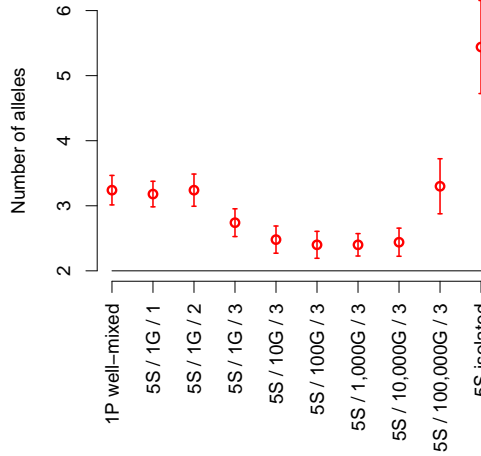
Figure 6.1: Number of alleles over time for different population sizes and well-mixed as well as structured populations with different migration rates.

uals migrating. For example, 5S / 10G / 4 in figure 6.2b stands for 5 subpopulations where migration occurs every 10 generations, and $10^{-4} \cdot m$ (where m is the size of the subpopulation) individuals emigrating from each subpopulation to any other connected subpopulation. This means that from subpopulation A , $10^{-4} \cdot m = 10^{-4} \cdot 10000/5 = 2$ individuals move to subpopulation B , a further 2 individuals emigrate to subpopulation C etc.

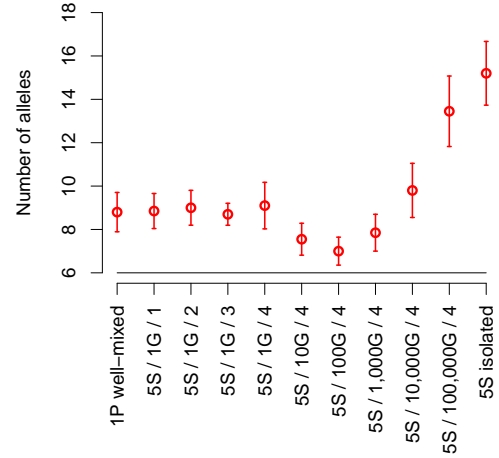
For total population sizes of 10000 (figure 6.2a) and 100000 (figure 6.2b), a minimum can be seen when plotting the number of alleles at the end of the simulation against different rates of migration. High migration rates maintain similar numbers of alleles when compared to the well-mixed population, but low migration rates maintain fewer alleles. Very low migration rates however can maintain even higher numbers of alleles, and a metapopulation where there is no migration between the subpopulations maintains the maximum number of alleles.

The situation is different for higher population sizes, as shown in figures 6.2c and 6.2d. For a total population size of 1 million (figure 6.2c), no drop in allele numbers can be seen for low levels of migration, but the same increase in allele numbers occurs for very low migration rates or no migration at all. Finally, for a total population size of 10 million (figure 6.2d), there is relatively little variation in the number of alleles at the end of the simulation span for all migration rates.

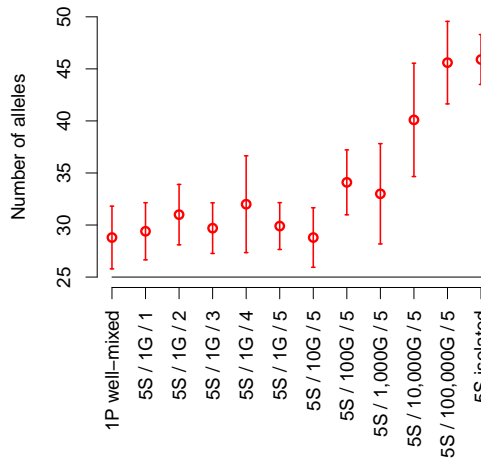
When looking at the trend for increasing numbers of subpopulations, only a slight increase can be observed for a migration rate of 2 individuals every generation (figure 6.3a), but if the migration interval is decreased to the minimum interval tested (100000 generations), then the number of alleles increases substantially when the metapopulation is subdivided into more subpopulations (figure 6.3b). The same is true if no migration between the subpopulations



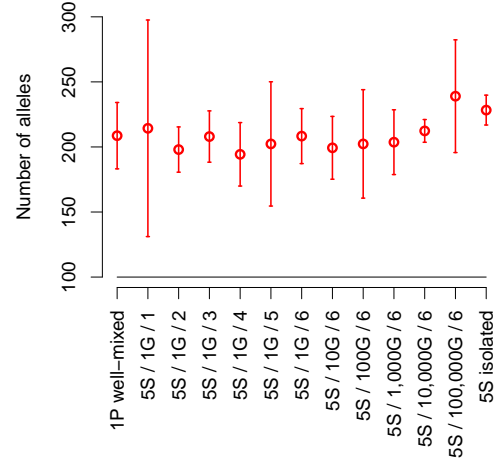
(a) 10000 ind. (total), 5 subpopulations



(b) 100000 ind. (total), 5 subpopulations



(c) 1 million ind. (total), 5 subpopulations



(d) 10 million ind. (total), 5 subpopulations

Figure 6.2: Number of alleles at the end of the simulation span for a metapopulation of 5 subpopulations in a star arrangement and different migration rates and population sizes.

occurs (figure D.1a), whereas no clear trend can be observed for most other migration rates (figure D.1b shows the number of alleles for a migration rate of 2 individuals every 100 generations).

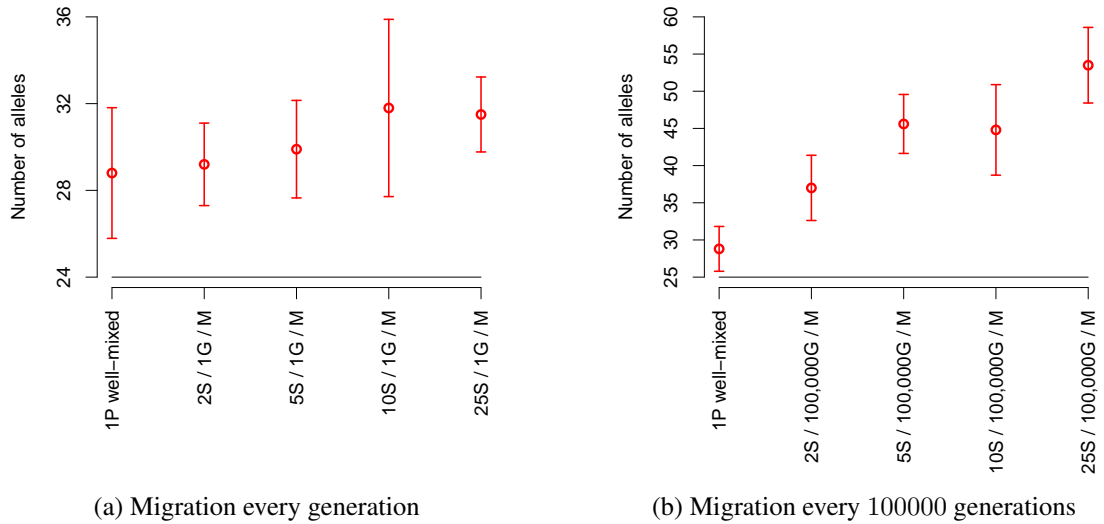


Figure 6.3: Number of alleles at the end of the simulation span for different numbers of subpopulations in a star arrangement and different migration rates. The total population size was 1 million.

6.3.2 Distribution of intrinsic merit of alleles

The distribution of intrinsic merit of alleles at the end of the simulation span was one of the major differences between the DAA model and traditional asymmetric overdominance models. For well-mixed populations, this distribution was found to be much wider for the DAA model, and a pronounced tail of alleles with low intrinsic merit existed. This section explores to what extent these observations apply to structured populations. For simplicity, the simple asymmetric overdominance model will subsequently be used as a representative of traditional asymmetric overdominance models, and will therefore be referred to as “asymmetric overdominance” or the “asymmetric overdominance model” only.

Figure 6.4 shows that different amounts of migration between subpopulations did not change the qualitative results from the well-mixed population. The DAA model still resulted in a much wider range of allele intrinsic merit compared to the asymmetric overdominance model, and additionally preserved the feature of a pronounced tail of alleles with low intrinsic merit. This result was robust across all population sizes and migration rates tested, and for both the star arrangement and the linear arrangement (for the latter, only the DAA model was tested), as figures D.6, D.7, D.8, D.9, D.10, D.11, D.12 and D.13 demonstrate.

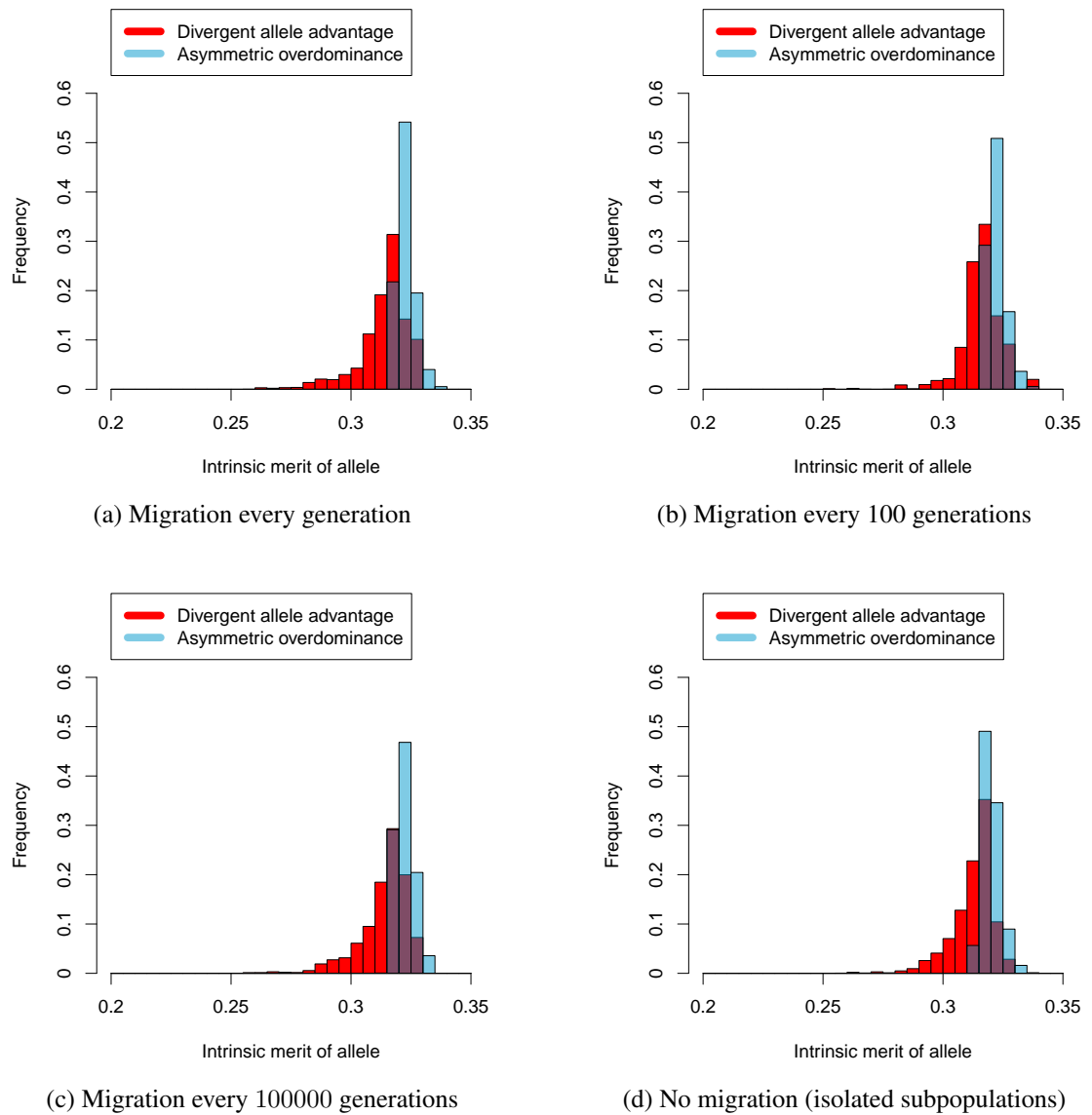


Figure 6.4: Distribution of intrinsic merit of alleles – comparison between the DAA model and the asymmetric overdominance model for a metapopulation of 5 subpopulations in a star arrangement, different migration rates and a population size of 100000.

6.3.2.1 Range of allele intrinsic merits

The range of allele intrinsic merits is a simple measure of diversity derived from the intrinsic merits of the final set of alleles already discussed and applied in section 5.3.1.4. Figures D.2 and D.3 show the allele intrinsic merit range for a metapopulation under the DAA model and different population sizes, migration intensities and numbers of subpopulations. As the DAA model was not very sensitive to the removal of transient alleles (section 5.3.1.4), all alleles were considered when calculating the intrinsic merit range.

Similar to the result for the total number of alleles, a minimum existed in the range of allele intrinsic merits at the end of the simulation versus different rates of migration for population sizes of 10000 (figure D.2a) and 100000 (figure D.2b). For higher population sizes, there was no strong dependency of the range of allele intrinsic merits on the different rates of migration (figures D.2c and D.2d).

The range of allele intrinsic merits further depended on the number of subpopulations that constitute the metapopulation, but while a fair amount of variability was present, the effect of population subdivision was rather weak and depended on the amount of migration between the subpopulations (figure D.3).

6.3.2.2 Weighted standard deviation in allele intrinsic merit

The weighted standard deviation in allele intrinsic merit responded similarly to the allele intrinsic merit range, both for the dependency on the migration between subpopulations (figure D.4) and on the number of subpopulations (figure D.5). A minimum existed for medium to small rates of migration for lower population sizes (figures D.4a and D.4b) while no strong dependency could be seen for higher population sizes (figures D.4c and D.4d). When varying the number of subpopulations, a slight to medium decrease in the weighted standard deviation in allele intrinsic merit towards a larger number of subpopulations could be observed for some medium to small rates of migration. As the variability of the weighted standard deviation in allele intrinsic merit was lower compared to the variability of the range of allele intrinsic merits, effects become more visible. Interestingly, the decrease in the weighted standard deviation for a larger number of subpopulations was most pronounced for an intermediate migration rate (figure D.5b), and less pronounced for a high migration rate (figure D.5a) or a very low migration rate (figure D.5c).

6.3.3 Diversity profiles

6.3.3.1 Naive diversity

As demonstrated in chapters 3, 4 and 5, more insights can be gained when not only single measures of diversity are compared, but entire diversity profiles. Figure 6.5 shows diversity profiles ${}^qD^I$ (Leinster and Cobbold, 2012). For $q = 0$ (which corresponds to the number of alleles), allele numbers were higher for very low levels of migration or metapopulations of isolated subpopulations, as figure 6.2 shows more explicitly. However, when moving towards $q = 2$, which is the inverse of expected homozygosity, then the differences between the migration rates became less pronounced, and for even higher values of q these differences became minimal. This was also true for lower population sizes, as shown in figures 6.5a and 6.5b. For a population size of 1 million (figure 6.5c), the diversity profiles can be separated into metapopulations consisting of isolated or nearly isolated subpopulations that exhibited higher diversity for $q \leq 2$ (green curves) and all other migration rates between subpopulations (red curves); within any of these two groups, differences were minor. When moving to a population size of 10 million, differences between migration rates increasingly disappeared (figure 6.5d).

When comparing metapopulations with the same total population size, but different numbers of subpopulations (this was only done for a total population size of 1 million), then a trend towards higher diversity for $q \leq 2$ can be seen when moving to larger numbers of subpopulations, particularly for metapopulations consisting of isolated or nearly isolated subpopulations (figure 6.6). The same figure also shows that for all levels of subdivision, metapopulations where migration is absent or nearly absent and metapopulations with more frequent migration separated into two groups, one of which was more diverse for low values of q .

Finally a comparison between the different arrangement types was performed, the star arrangement with maximum connectivity and the linear arrangement with low connectivity. While in many cases differences between these arrangements were minor (figure D.14 shows population sizes of 1 million and 10 millions and two different migration rates, while figure D.15 shows all population sizes tested and a migration rate of 2 individuals every 100 generations), there was nevertheless an interesting dependency on the migration interval for lower population sizes. Figure 6.7 demonstrates that diversity for lower values of q was higher in the star arrangement when migration is frequent, but when migration is extremely infrequent, then diversity was higher in the linear arrangement.

Nevertheless, while small differences between different rates of migrations existed, these were dwarfed by differences between overdominance models themselves (figures D.16, D.17, D.18 and D.19), which were much more significant for all amounts of migration tested.

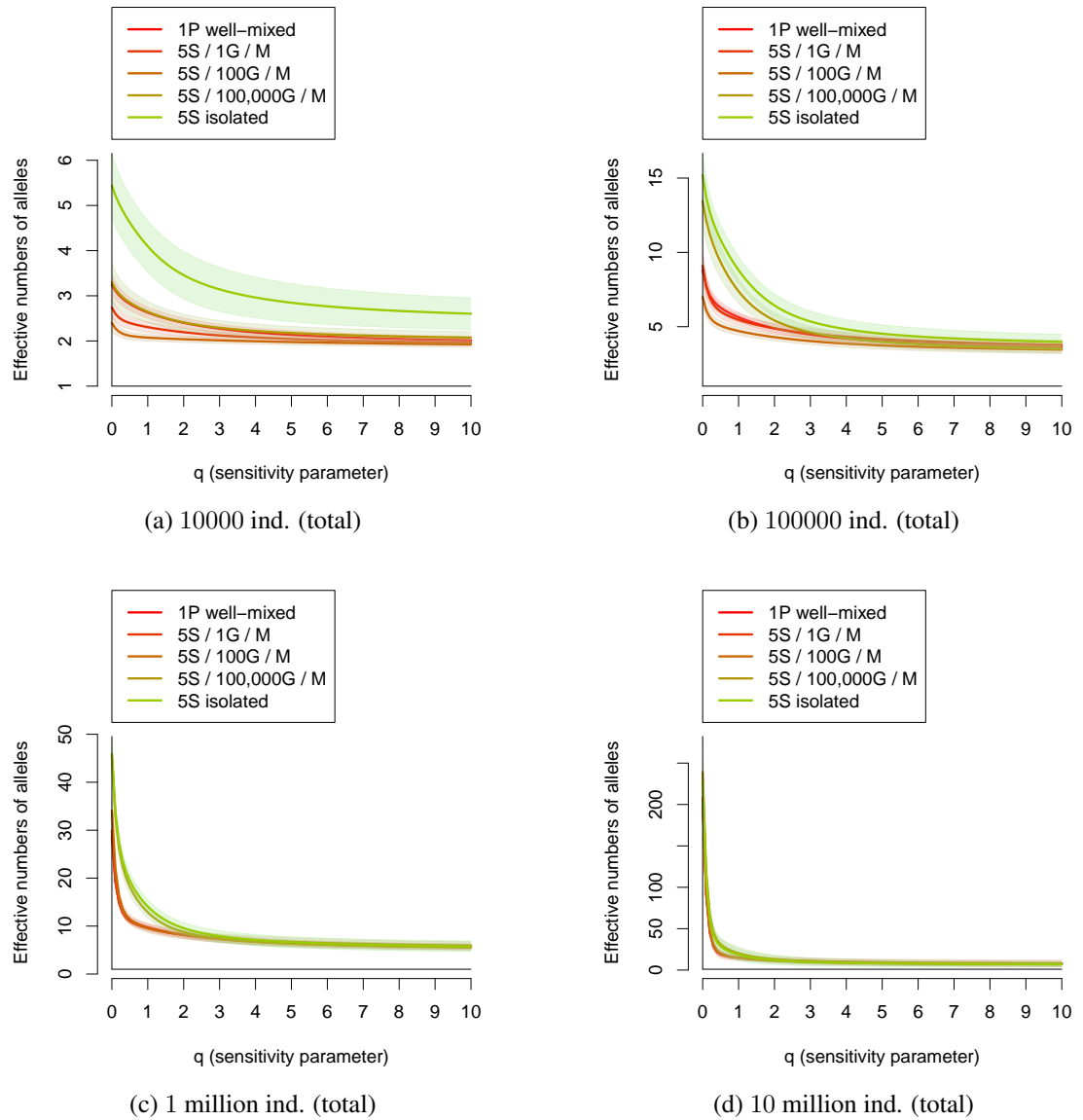
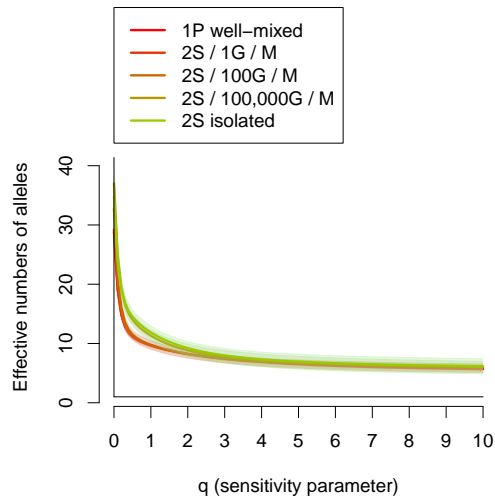
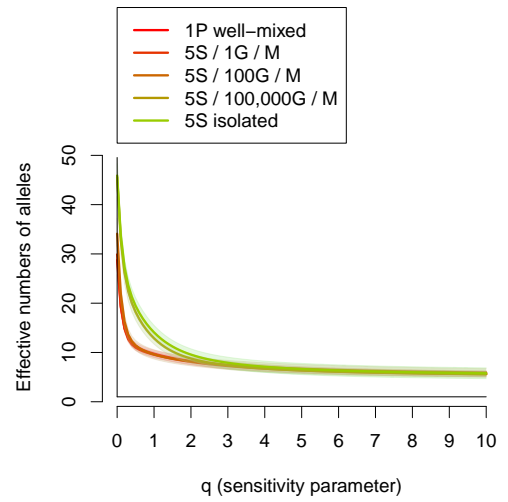


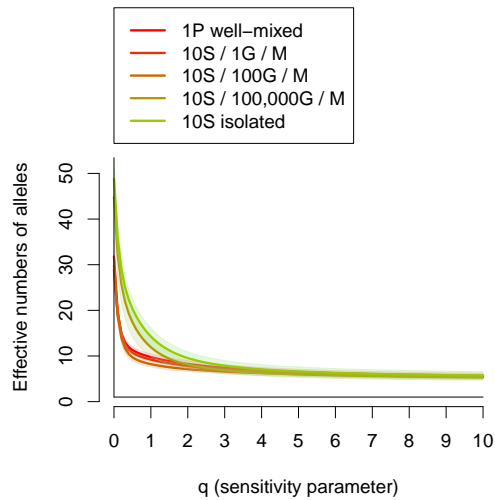
Figure 6.5: Diversity profiles (naive diversity) – comparison between different migration rates for a metapopulation of 5 subpopulations in a star arrangement and all population sizes tested.



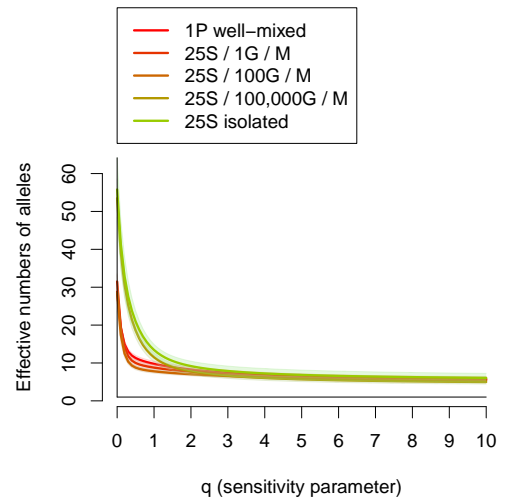
(a) 1 million ind. (total), 2 subpopulations



(b) 1 million ind. (total), 5 subpopulations



(c) 1 million ind. (total), 10 subpopulations



(d) 1 million ind. (total), 25 subpopulations

Figure 6.6: Diversity profiles (naive diversity) – comparison between metapopulations consisting of different numbers of subpopulations in a star arrangement, and different migration intervals.

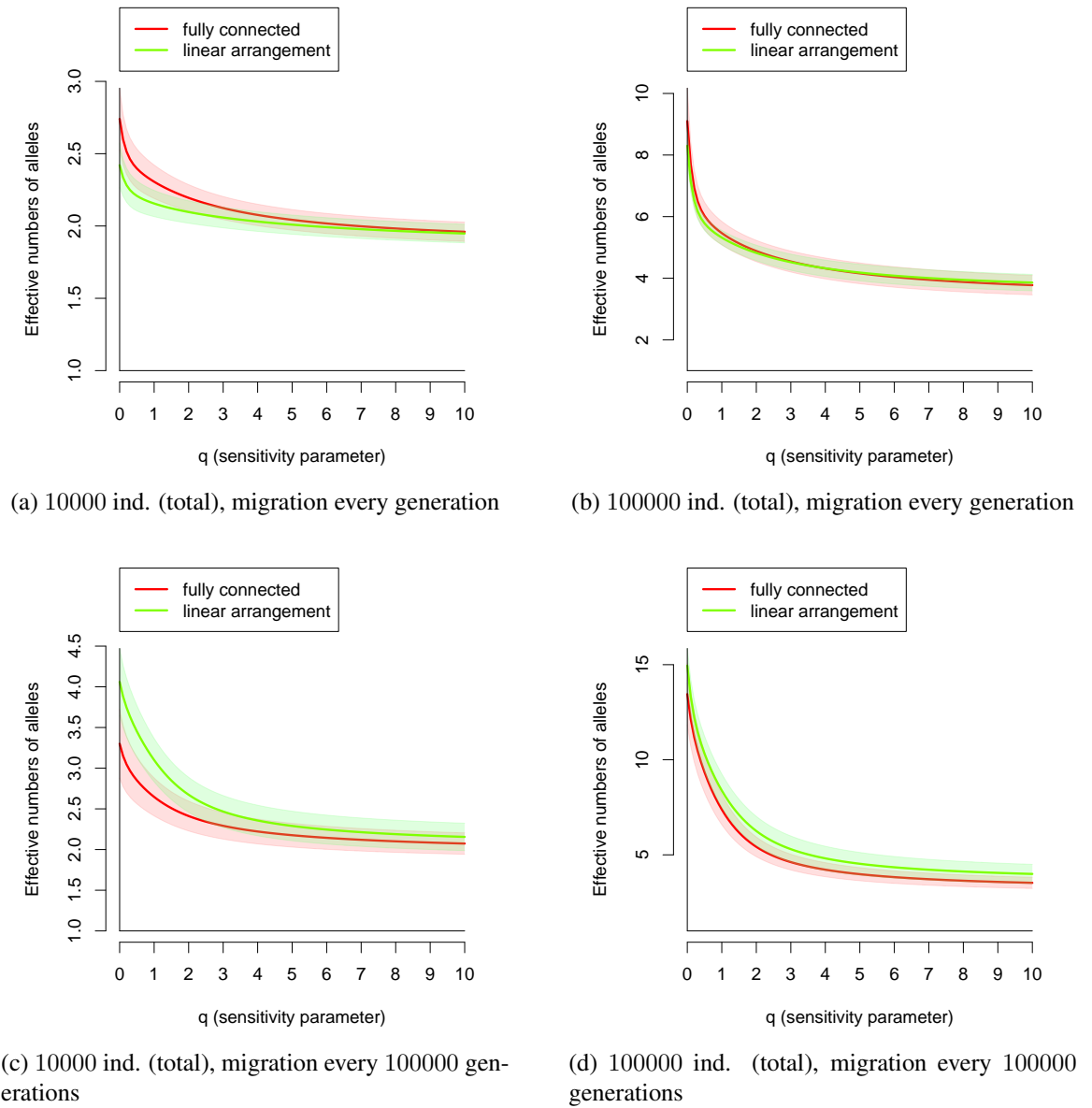


Figure 6.7: Diversity profiles (naive diversity) – comparison between a metapopulation of 5 subpopulations in a star arrangement and a metapopulation of 5 subpopulations in a linear arrangement, for a population size of 10000 (left) and 100000 (right) respectively.

6.3.3.2 Similarity-sensitive diversity

In the case of the well-mixed population, the higher variation in intrinsic merit for the DAA model resulted in a higher similarity-sensitive diversity (with differences in intrinsic merit as the similarity measure) for this model compared to the asymmetric overdominance model, even though its naive diversity was lower. This was an important result, as it demonstrated that the DAA model was capable of both maintaining large numbers of alleles and variation among them, expressed by the range of intrinsic merits of the extant alleles. This section examines whether the same result holds for structured populations.

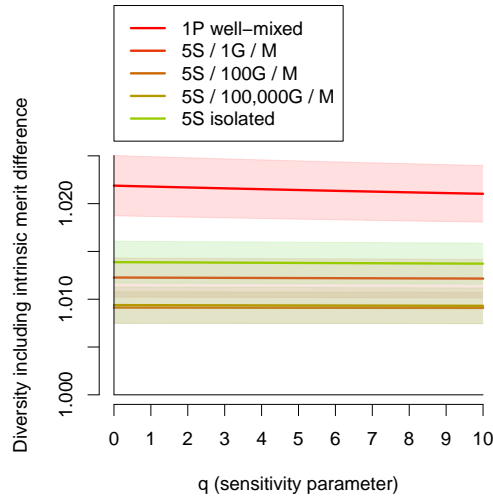
Section 6.3.2 already demonstrated that the distributions of the intrinsic merits of alleles are similar across migration rates. Now the dependency of similarity-sensitive diversity based on the difference between intrinsic merits of alleles on the migration rate is explored. Figure 6.8 shows how this similarity-sensitive diversity varied across migration intervals. Differences were more pronounced when either the population size was low (figure 6.8a) or the number of subpopulations in the metapopulation was high (figure 6.8d). In both cases, diversity was highest for the well-mixed population and lowest for a metapopulation with infrequent but existing migration. When comparing the DAA model with the asymmetric overdominance model, then this similarity-sensitive diversity was substantially and consistently higher in the former (figures D.20, D.21, D.22 and D.23), just as in the well-mixed scenario.

Finally, diversity profiles that include the number of differences between the sequences of the two alleles on the considered locus (sequence divergence, figure 6.9), were compared while the migration interval was varied. While naive diversity was highest for the lowest levels of migration or metapopulations of isolated subpopulations, diversity that accounts for sequence divergence was highest for the well-mixed single population and the population with the highest level of migration.

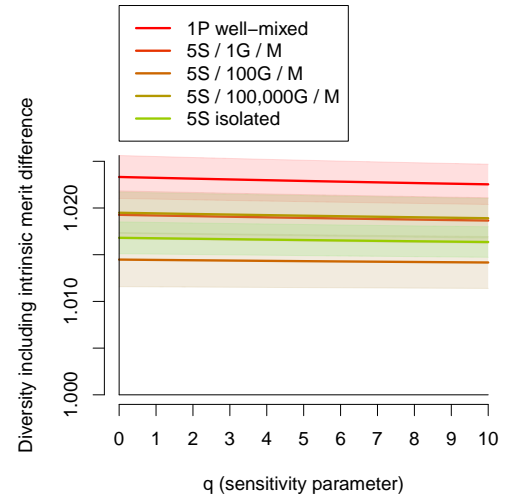
Comparing the DAA model with the asymmetric overdominance model in terms of diversity based on sequence divergence yielded a qualitatively similar result to the well-mixed scenario, with a much higher similarity-sensitive diversity for the DAA model as compared to the asymmetric overdominance model, particularly for large population sizes (figures D.24, D.25, D.26 and D.27).

6.3.3.3 Comparison between model results and observed allele frequencies

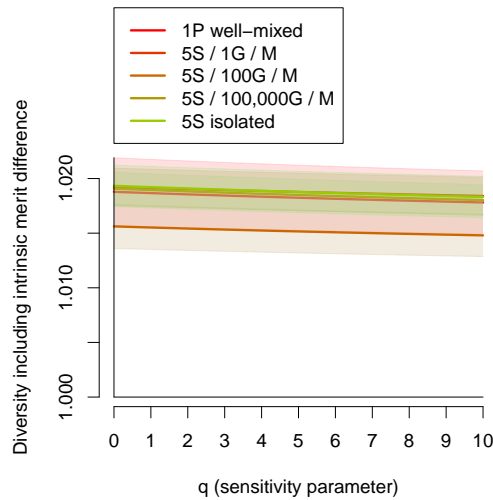
As in the well-mixed scenarios, the distribution of allele frequencies was compared between natural populations (table 4.3) and simulation results using various overdominance models. Diversity measures applied were the expected heterozygosity, the share of the 5 most frequent alleles and diversity profiles only using the allele frequencies (i.e. profiles of naive diversity).



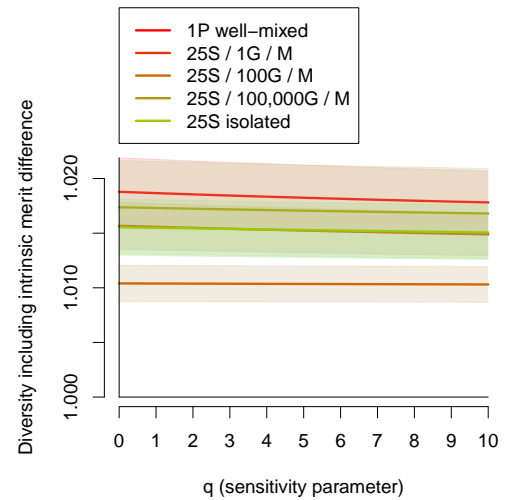
(a) 10000 ind. (total), 5 subpopulations (star)



(b) 100000 ind. (total), 5 subpopulations (star)

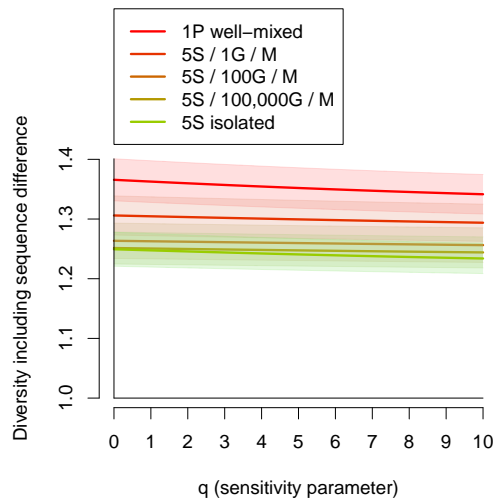


(c) 1 million ind. (total), 5 subpopulations (star)

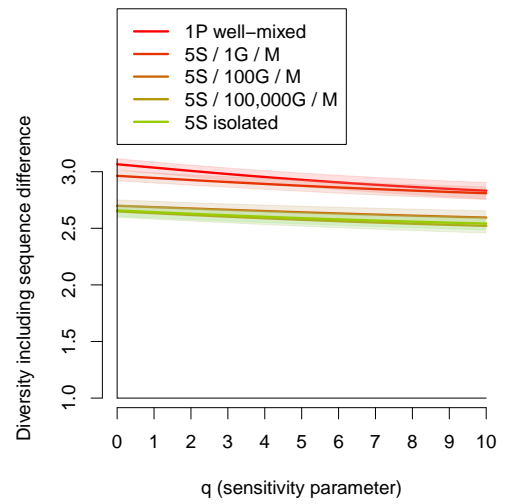


(d) 1 million ind. (total), 25 subpopulations (star)

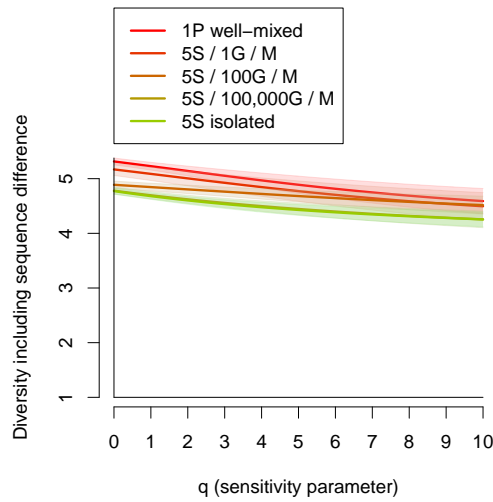
Figure 6.8: Diversity profiles (diversity that accounts for intrinsic merit differences between the alleles) – comparison between different migration rates, population sizes and numbers of subpopulations.



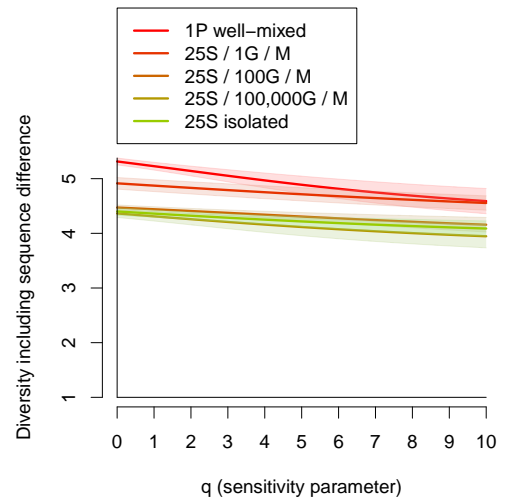
(a) 10000 ind. (total), 5 subpopulations (star)



(b) 100000 ind. (total), 5 subpopulations (star)



(c) 1 million ind. (total), 5 subpopulations (star)



(d) 1 million ind. (total), 25 subpopulations (star)

Figure 6.9: Diversity profiles (diversity that accounts for allele sequence divergence) – comparison between different migration rates, population sizes and numbers of subpopulations.

The first diversity measure applied was the expected heterozygosity. Figure 6.10 compares expected heterozygosities of the DAA model and the asymmetric overdominance model with observational data (see table 4.3). Similar to the well-mixed scenario for this setting, the DAA model predicted observed values much better than the asymmetric overdominance model, which overestimated heterozygosity particularly for a metapopulation of 5 isolated subpopulations (figure 6.10d). As mentioned in the methods section (6.2.3), expected heterozygosities were calculated according to the first approach mentioned (i.e. the gene pools of the subpopulations were merged before calculating the expected heterozygosities).

The connectivity of the subpopulations in a structured population, however, did not result in noticeable differences between the expected heterozygosities, with both the star arrangement and the linear arrangement predicting expected heterozygosities of natural populations well (figure D.28), and differences between the arrangements were not greater than the stochastic variation.

Next the share of the 5 most frequent alleles of natural populations was compared to the simulation output of the DAA model and the asymmetric overdominance model (figure 6.11). The comparison again shows that sensitivity towards the structure of the population and the amount of migration is small compared to the influence of the overdominance model used. Also in terms of the share of the 5 most frequent alleles, the DAA model was much closer to observed values than the asymmetric overdominance model, with predictions of the latter being particularly poor for a metapopulation of 5 isolated subpopulations (figure 6.11d). Just as in case of expected heterozygosities, the cumulative share of the 5 most frequent alleles was calculated from the merged gene pools of the subpopulations.

The influence of the arrangement of the subpopulations on the share of the 5 most frequent alleles was weak, with both arrangement types giving virtually the same results (figure D.29).

Finally, diversity profiles $^qD^I$ calculated with samples from natural populations were compared to diversity profiles of simulation results, obtained from allele frequencies of the final gene pool. First, the DAA model was compared to the asymmetric overdominance model for a metapopulation of 5 well-connected subpopulations and different migration rates (figure 6.12). While the highest migration rate (figure 6.12a) gave a similar result to the well-mixed situation, an even better fit of the DAA model predictions was obtained for lower migration rates (figures 6.12b and 6.12c) or no migration (figure 6.12d), while the asymmetric overdominance model did not fit the observed allele frequencies for any migration rate tested.

As was the case for expected heterozygosities and the share of the 5 most frequent alleles, there was little variation in the results if the star arrangement was changed to the linear arrangement, with both giving a similar fit to the naive diversities calculated from observed allele frequencies (figure D.30), and differences between the arrangements were in the order of the intrinsic stochastic variation.

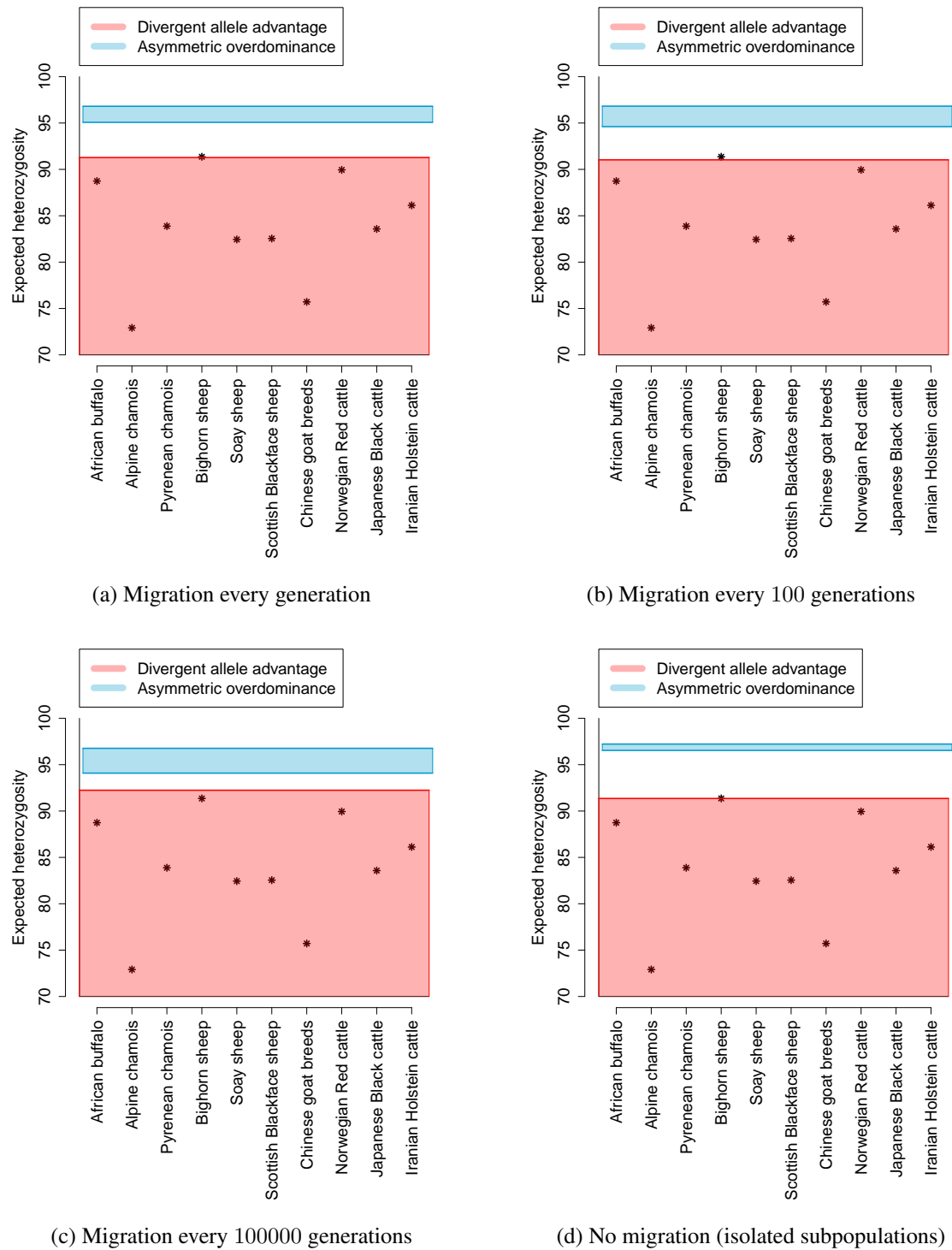


Figure 6.10: Expected heterozygosity for predictions from different overdominance models and observed data for a metapopulation of 5 subpopulations in a star arrangement and various migration rates. The red and blue bands span population sizes from 10000 to 10 million.

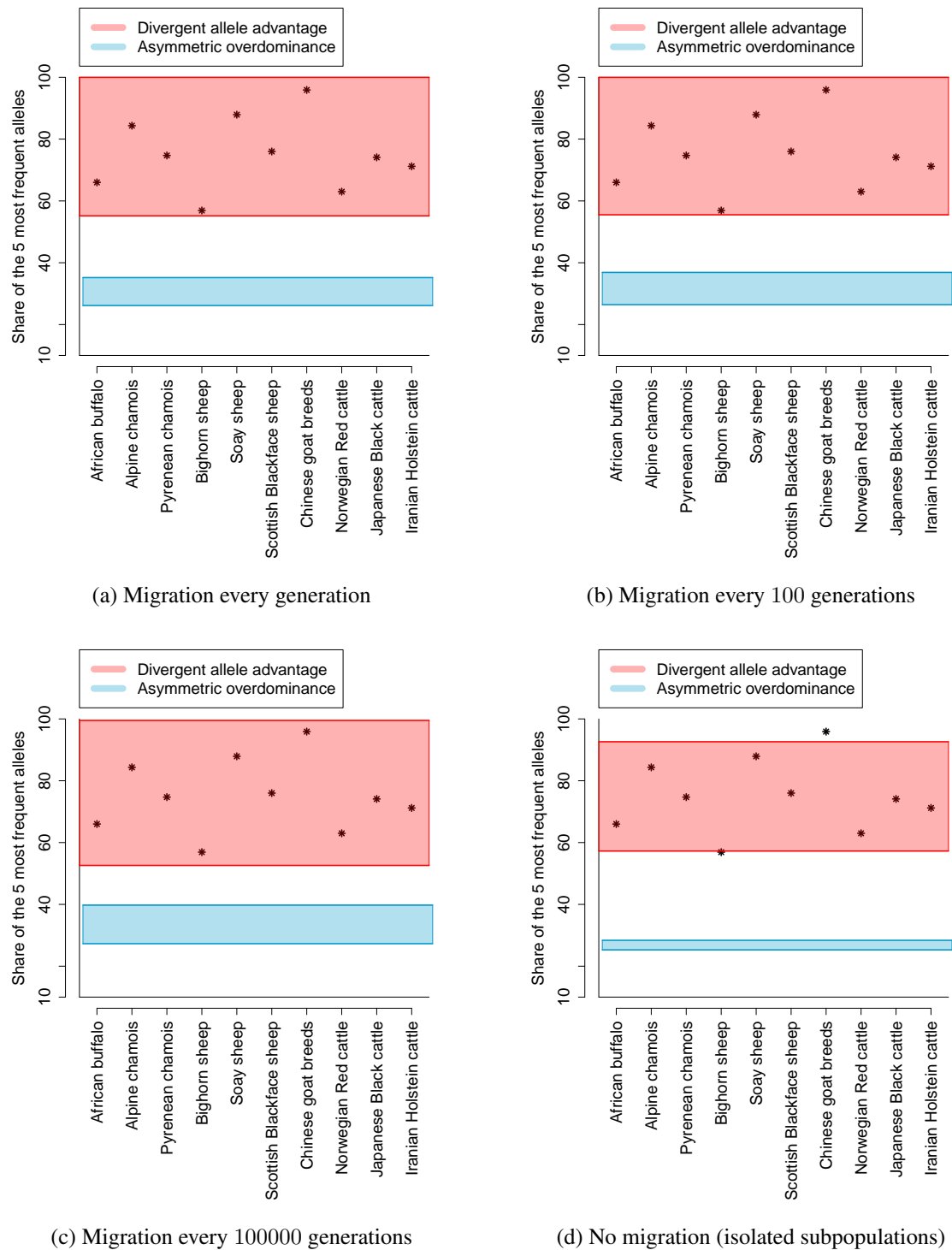


Figure 6.11: Share of the five most frequent alleles for predictions from different overdominance models and observed data for a metapopulation of 5 subpopulations in a star arrangement and various migration rates. The red and blue bands span population sizes from 10000 to 10 million.

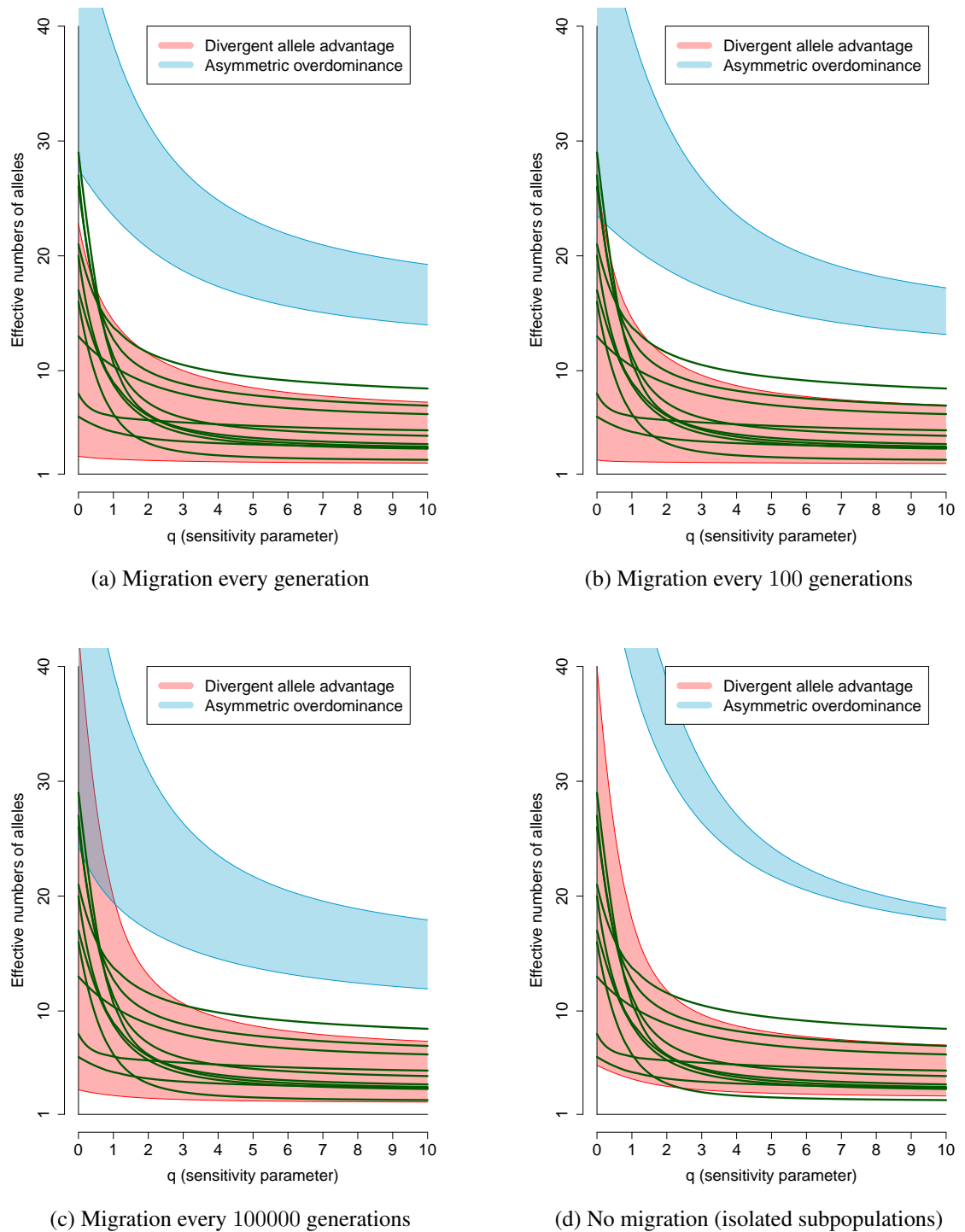


Figure 6.12: Diversity profiles (naive diversity) for predictions from different overdominance models and observed data for a metapopulation of 5 subpopulations in a star arrangement and various migration rates. The red and blue bands span population sizes from 10000 to 10 million.

6.3.4 Trans-species polymorphism

In single, well-mixed populations, the final gene pool of simulations using the DAA model exhibited a mix of alleles that emerged early and alleles with more recent emergence times (section 4.4.2.7); this feature was considered consistent with trans-species polymorphism.

For structured populations, allele emergence times were recorded for the DAA model for both the star (figures D.31, D.32, D.33 and D.34) and the linear arrangement (figures D.35, D.36, D.37 and D.38), and for various migration rates. These figures demonstrate that again alleles spanning the whole range of emergence times existed in the gene pool at the end of the simulation, just as in the single, well-mixed population (section 4.4.2.7). Therefore, the DAA model is also consistent with trans-species polymorphism when the population is structured.

However, one notable peculiarity of the emergence times of a structured population as compared to the well-mixed scenario stood out. Particularly for larger total population sizes (i.e. population sizes of 1 million and 10 million), a pronounced peak in the allele emergence time interval of 5 million to 6 million generations, just after the split of the single population into subpopulations, could be observed. This is true for both the star and the linear arrangement, and for all migration rates tested (figures D.31, D.32, D.33, D.34, D.35, D.36, D.37 and D.38). It appears that just after the gene pool of the single population had been subdivided into subpopulations, there were favourable conditions for newly mutated alleles to invade the gene pool of the respective subpopulation. As the main aim of this chapter is to confirm that the results obtained for the well-mixed population are also valid when the population is structured, the reason underlying this observation has not been investigated.

6.4 Discussion and Conclusions

Many populations show a certain degree of spatial structure, with barriers limiting migration between different parts of the population. To validate the conclusions of chapter 4, a simple metapopulation model that accounts for population structure was investigated.

A number of assumptions were made in the metapopulation model that simplify the structure of real-world populations. In the simulations, a single population was divided into subpopulations that were themselves well-mixed, all subpopulations emerged at the same time when the original population was split, all these subpopulations were of the same size, and migration between any two subpopulations was symmetric in terms of the number of individuals migrating. Despite these simplifications, the model nevertheless captured another important feature of real populations, the migration between subpopulations, and thus this expanded model accounts for all major evolutionary forces.

The results demonstrate that there is a nonlinear response in the number of alleles maintained when migration between subpopulations is varied from high to low, with the minimum number of alleles occurring for low but not very low migration rates in populations up to a total population size of 100000. In these scenarios with low levels of migration every 10 to 1000 generations, genetic drift is relatively strong due to the small sizes of the subpopulations (2000 individuals and 20000 individuals for total population sizes of 10000 and 100000, respectively), but migration is still strong enough to prevent independent evolution of the subpopulations.

When there is migration between subpopulations, even if the migration rate is very low, then subpopulations barely diverge. The effect of drift however is strong when the population size is low, therefore allele numbers in a metapopulation consisting of small subpopulations are negatively affected by a very low amount of migration. For larger population sizes of 1 million and above, however, a metapopulation of 5 subpopulations does not exhibit a significantly lower number of alleles compared to a single, well-mixed population of the same size, as genetic drift is weaker relative to the other evolutionary forces. The higher “carrying capacity” for alleles due to weaker competition between them in the subpopulations (as compared to the single population) and genetic drift are more in balance, which results in similar numbers of alleles being supported.

It is indeed important to keep in mind that the number of alleles maintained is a result of interactions between all evolutionary forces, so not only the amount of migration and genetic drift play a role, but also selection and mutation. All these forces acting in common lead to the observed dependency on population size.

Other diversity measures however show only little variation between the single, well-mixed population and the metapopulation, both in terms of migration rates and connectivity, and results that have been established in chapter 4 still hold in a metapopulation scenario likewise. This adds to the evidence that heterozygote advantage in the form of divergent allele advantage is the main driver of MHC diversity, and that the novel model based on divergent allele advantage is the best available model for explaining features associated with the MHC polymorphism.

Chapter 7

General Conclusions

This research has addressed the question of why genes encoding the MHC are the most polymorphic loci in vertebrates. Among various hypotheses, balancing selection mediated by pathogens appears to be the most promising explanation for the main force driving diversity at the vertebrate MHC. Heterozygote advantage is the only concept of the three principal explanations of pathogen-mediated balancing selection on MHC genes that has a clear mechanistic footing based on the function of the molecules, since heterozygotes can bind a broader spectrum of peptides than can homozygotes due to their broader immune surveillance.

Pursuing this approach leads to a special form of heterozygote advantage, called divergent allele advantage (Wakeland et al., 1990), whose core idea is that heterozygotes with divergent alleles recognise a wider set of peptides derived from parasites than heterozygotes with similar alleles, or even homozygotes; this implies more effective immune responses against each parasite and also against a wider range of parasites. Models based on this hypothesis were examined both at equilibrium, with allele frequency changes every generation due to selection, and in scenarios that additionally included mutation, drift, and even migration between subpopulations. These divergent allele advantage models were compared to traditional overdominance models and observed data, to check to what extent and under which conditions they can explain the main features of MHC polymorphism. A notably useful but so far infrequently utilised method, the diversity profile ${}^qD^Z$ (Leinster and Cobbold, 2012), helped to characterise MHC diversity at multiple scales and provided a valuable tool for comparisons.

Understanding how the high levels of MHC diversity are maintained is of significant practical value, as the MHC is the genetic region central to conferring resistance to pathogens in vertebrates (Pirotney and Oliver, 2006), so exploiting its potential could markedly improve population health of managed populations. The novel model based on the DAA hypothesis presented in chapters 4 and 5 suggests that one way to achieve this improvement is to use sequence data of MHC alleles and choose mating partners based on the compatibility of

their genotypes. This would result in offspring that has divergent alleles on one or more MHC loci, making these animals more resistant against parasites according to the DAA hypothesis. If additionally the fitnesses that alleles confer (the intrinsic merits of the alleles) can be estimated, then this information could be included in a selection index together with sequence-based information, which could potentially further improve the fitness of the offspring.

The results of the analysis indicate that divergent allele advantage has the potential to be a fundamental driver of MHC diversity. This special form of heterozygote advantage produced results that were consistent with all the major features of MHC diversity, including the large number of alleles, the intermediate level of homozygosity, or more general, the characteristic shape of the diversity profile, the variation in fitness contributions of alleles and the similarities among alleles from different species, i.e. the trans-species polymorphism. It is therefore not necessary to assume an essential role for other mechanisms that could potentially create or maintain diversity at the MHC, although it is likely that these mechanisms will additionally be present and will therefore modify diversity patterns.

The conclusion that the modelling presented here is consistent with a fundamental role for divergent allele advantage in the maintenance of MHC diversity could have – if empirically verified – important implications. While great progress has been made in selective breeding of farm animals for production traits, the potential to increase disease resistance could not be fully exploited, mainly because classic selection schemes do not deal with intra-locus interactions such as heterozygote advantage. As the MHC is central to the recognition of pathogens, being able to use this region could be the key to substantially improve disease resistance of managed animals, including wild animals of conservation concern.

Another implication is that the modelling predicts that a relatively large proportion of alleles (counting allele representatives rather than allele type) has a low intrinsic merit in the DAA model, which would make a significant part of any population relatively unfit. Identifying and managing, or, where necessary, even eliminating these alleles in breeding programmes could lead to a marked improvement in disease resistance of a population. Further, it was demonstrated that the recovery of population fitness can be a slow process under divergent allele advantage, given that allelic lineages were lost. This would indicate that there is a potential that a population of conservation concern is at risk of increased disease susceptibility if such loss of lineages occurs. Finally, the approach presented here may serve as a model for “post genomic” analysis of other loci.

These results together might substantially improve our understanding of selection at the MHC, and could enhance disease control of wild and domestic animal populations, but also management of human health, given that experimental evidence strengthens the case for the novel model based on the DAA hypothesis.

Appendix A

Additional figures for chapter 3

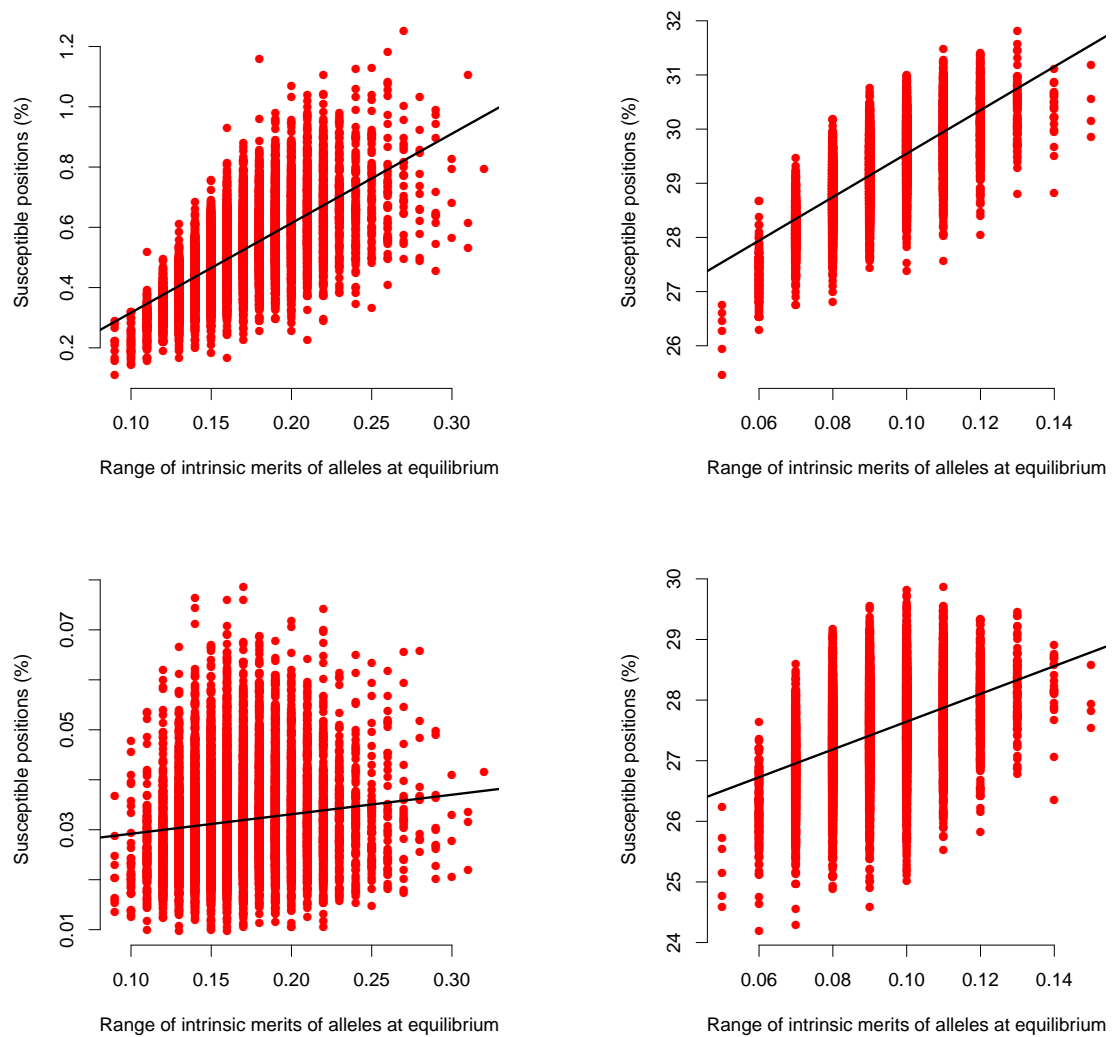


Figure A.1: Average proportion of susceptible recognition sites of heterozygotes vs. range of intrinsic merits of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). In the top row all persisting alleles are weighted equally, whereas in the bottom row the proportions of the heterozygous genotypes are used as weights. A positive correlation can be seen for all situations.

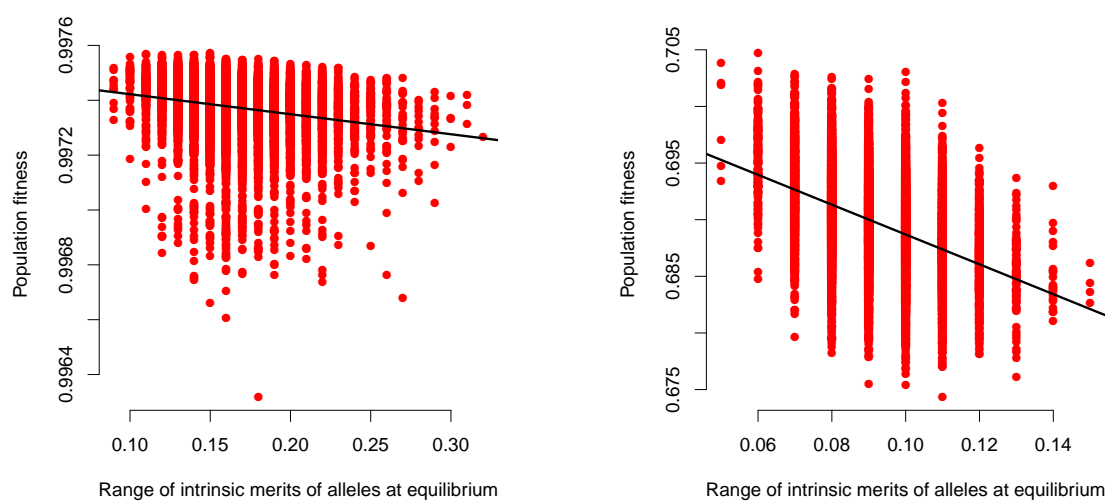
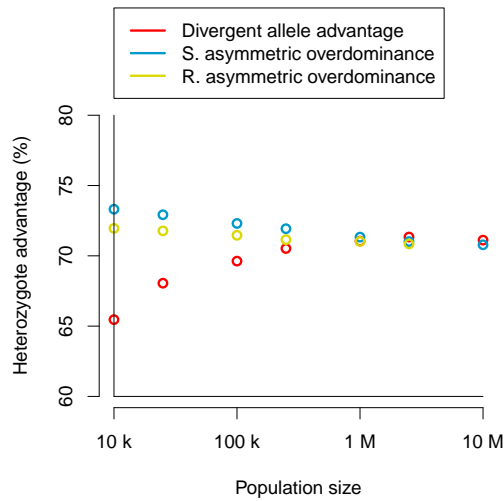


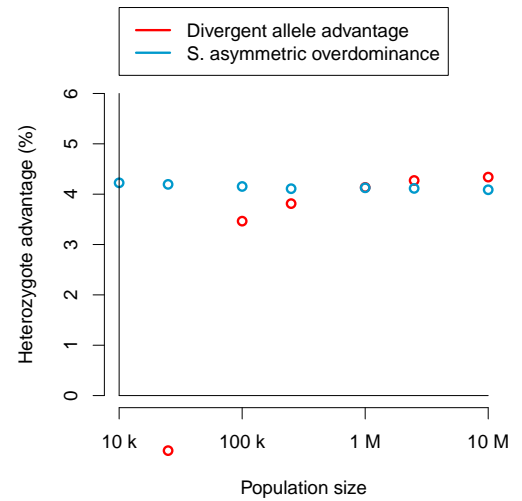
Figure A.2: Population fitness vs. range of intrinsic merits of alleles maintained at equilibrium for scenario 1 (left) and scenario 2 (right). A negative correlation can be seen for both scenarios.

Appendix B

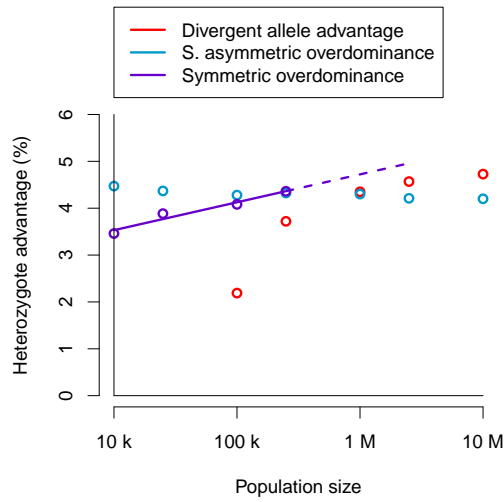
Additional figures for chapter 4



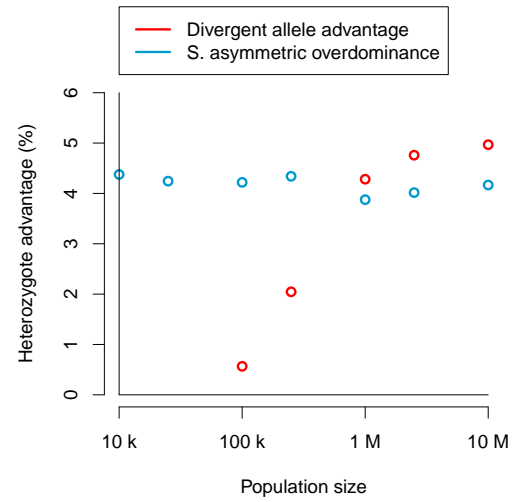
(a) Setting 3



(b) Setting 4



(c) Setting 5



(d) Setting 6

Figure B.1: Heterozygote advantage for different overdominance models and settings 3, 4, 5 and 6.

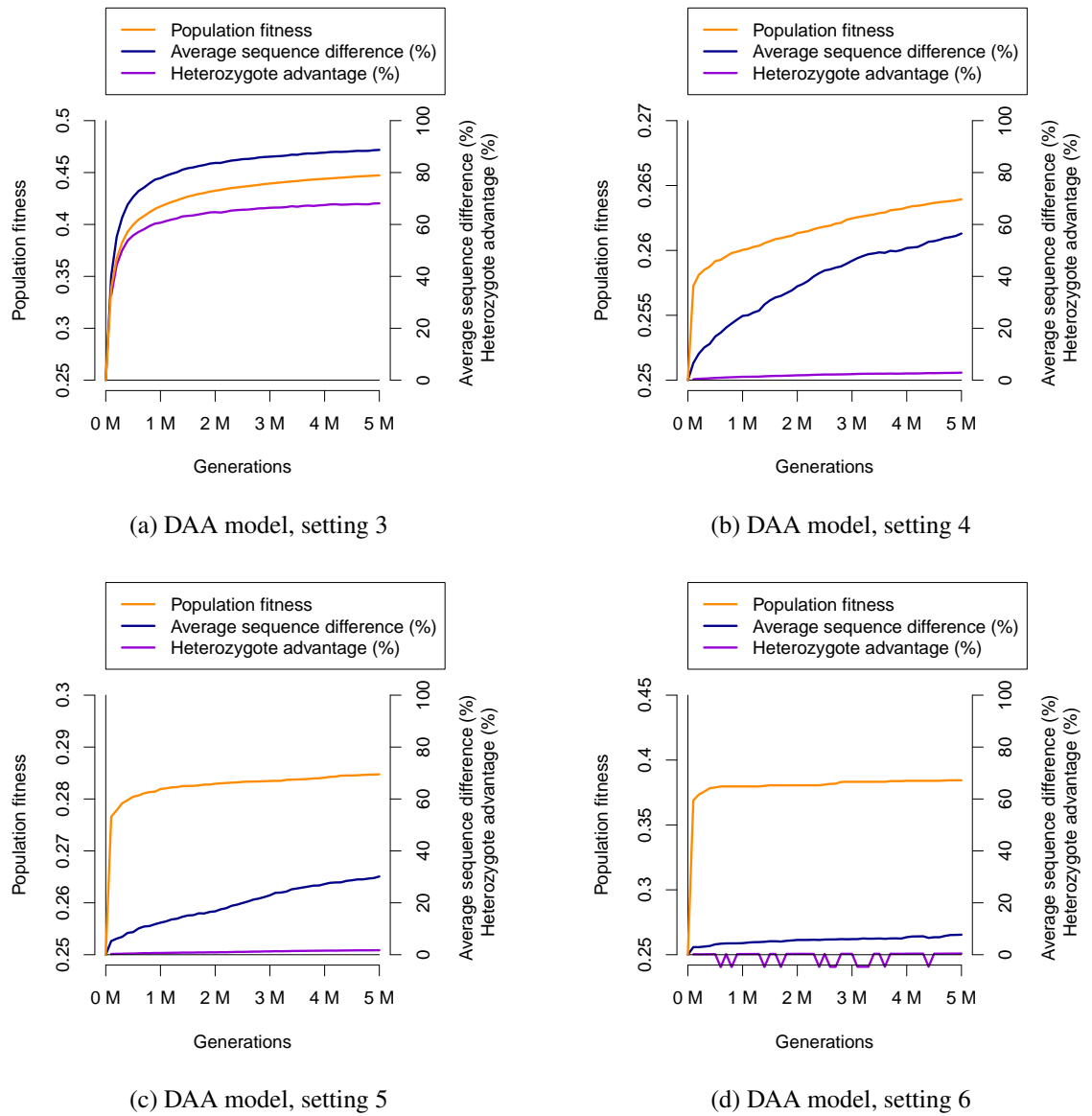


Figure B.2: Average sequence difference, heterozygote advantage and population fitness for the DAA model, settings 3, 4, 5 and 6 and a population size of 100000.

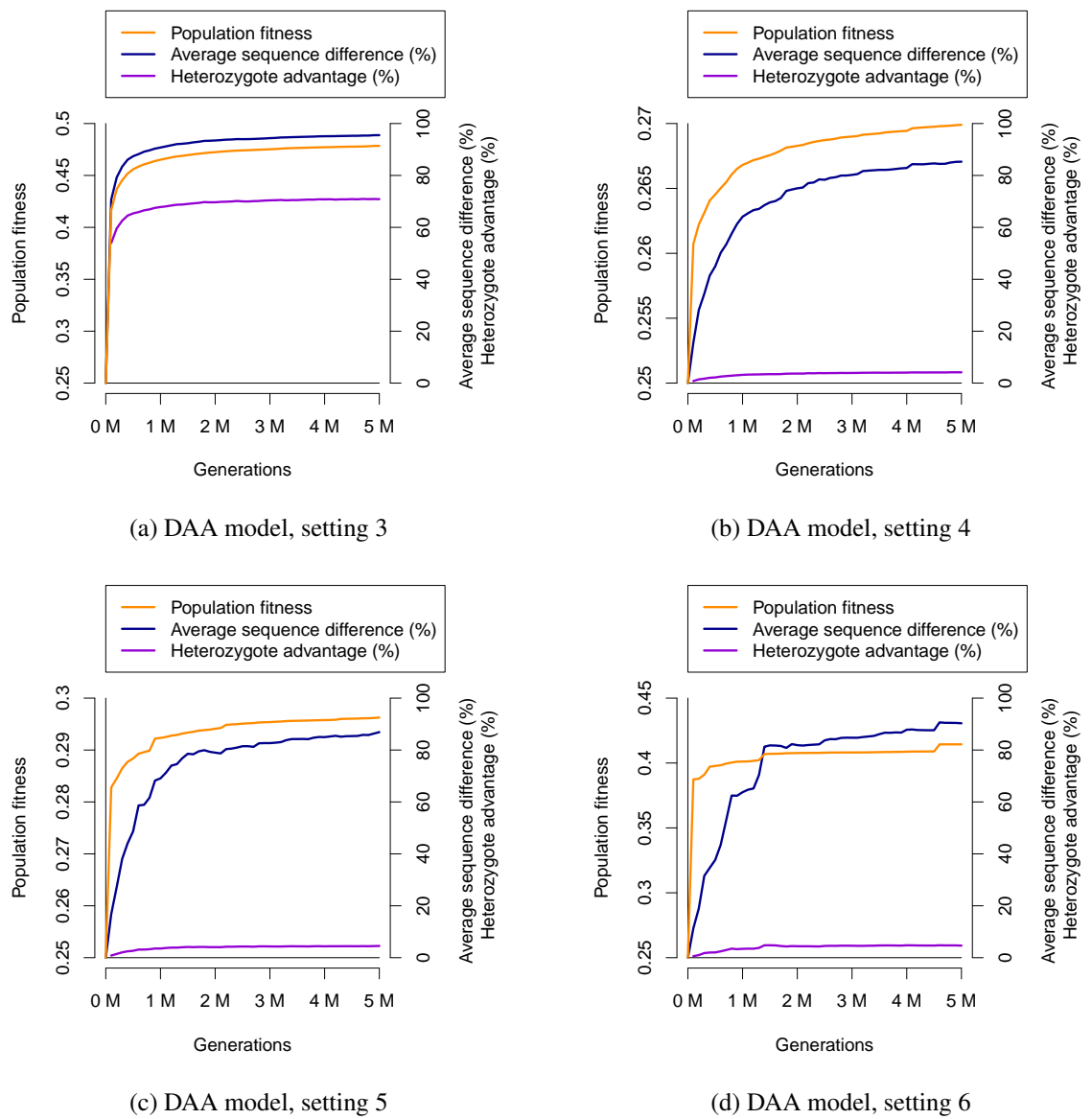
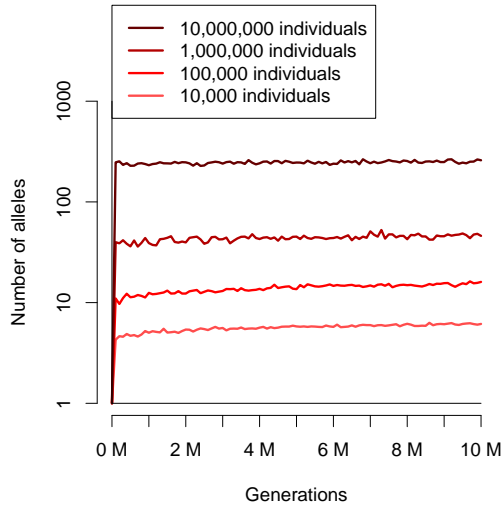
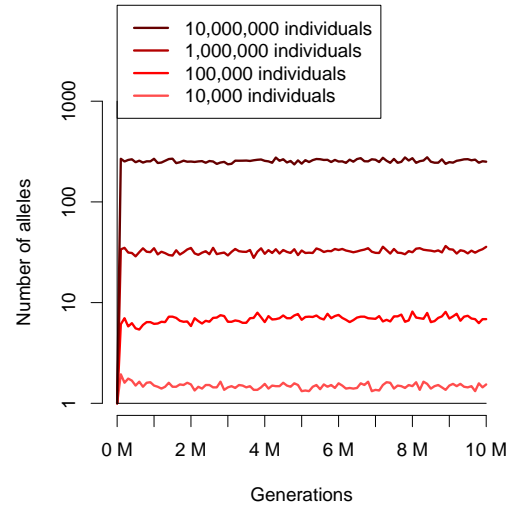


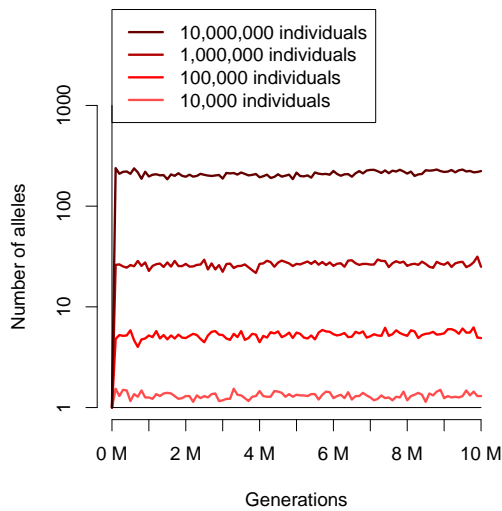
Figure B.3: Average sequence difference, heterozygote advantage and population fitness for the DAA model, settings 3, 4, 5 and 6 and a population size of 10 million.



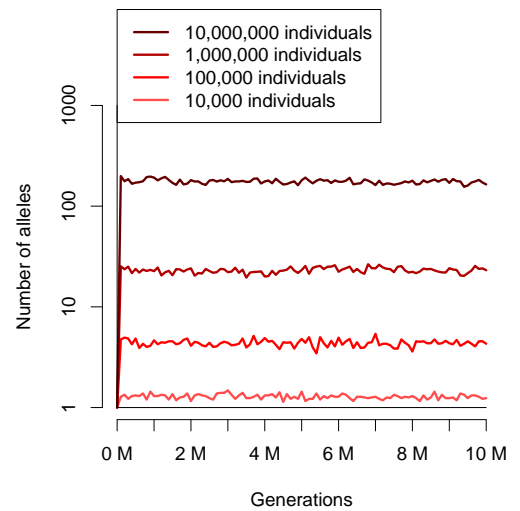
(a) Setting 3



(b) Setting 4

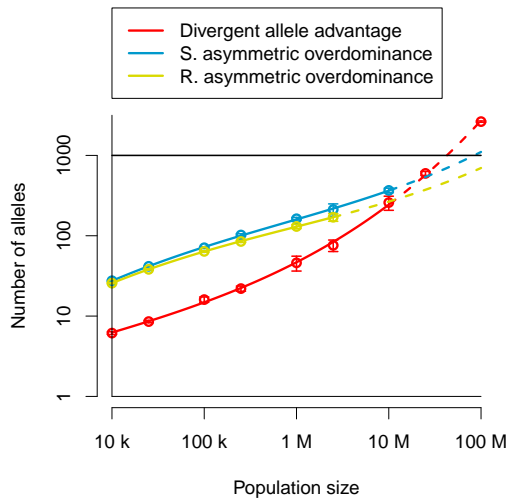


(c) Setting 5

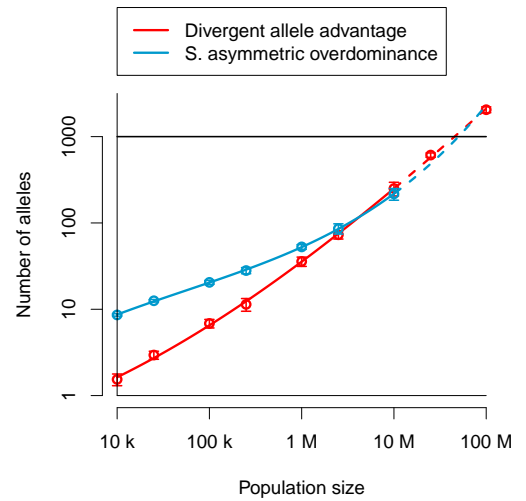


(d) Setting 6

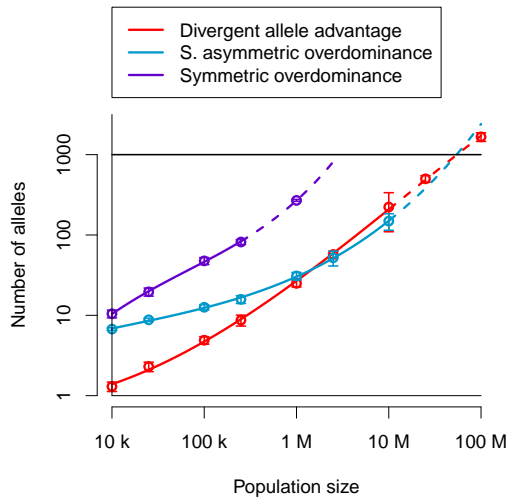
Figure B.4: Number of alleles over time in the DAA model for settings 3, 4, 5 and 6. Shades of red represent different sized populations. The simulation starts with a single allele and models the action of mutation, selection and drift over millions of generations.



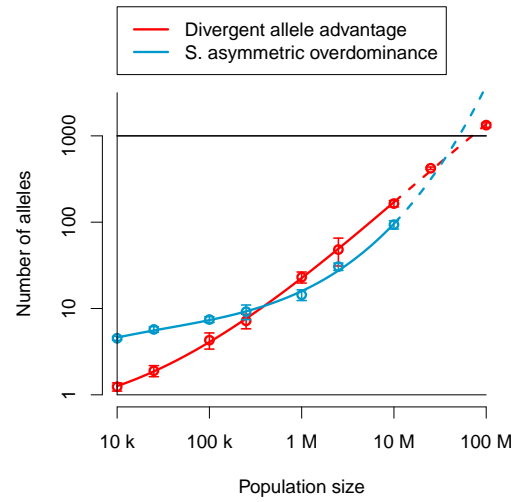
(a) Setting 3



(b) Setting 4



(c) Setting 5



(d) Setting 6

Figure B.5: Number of alleles vs. population size for settings 3, 4, 5 and 6 and different overdominance models.

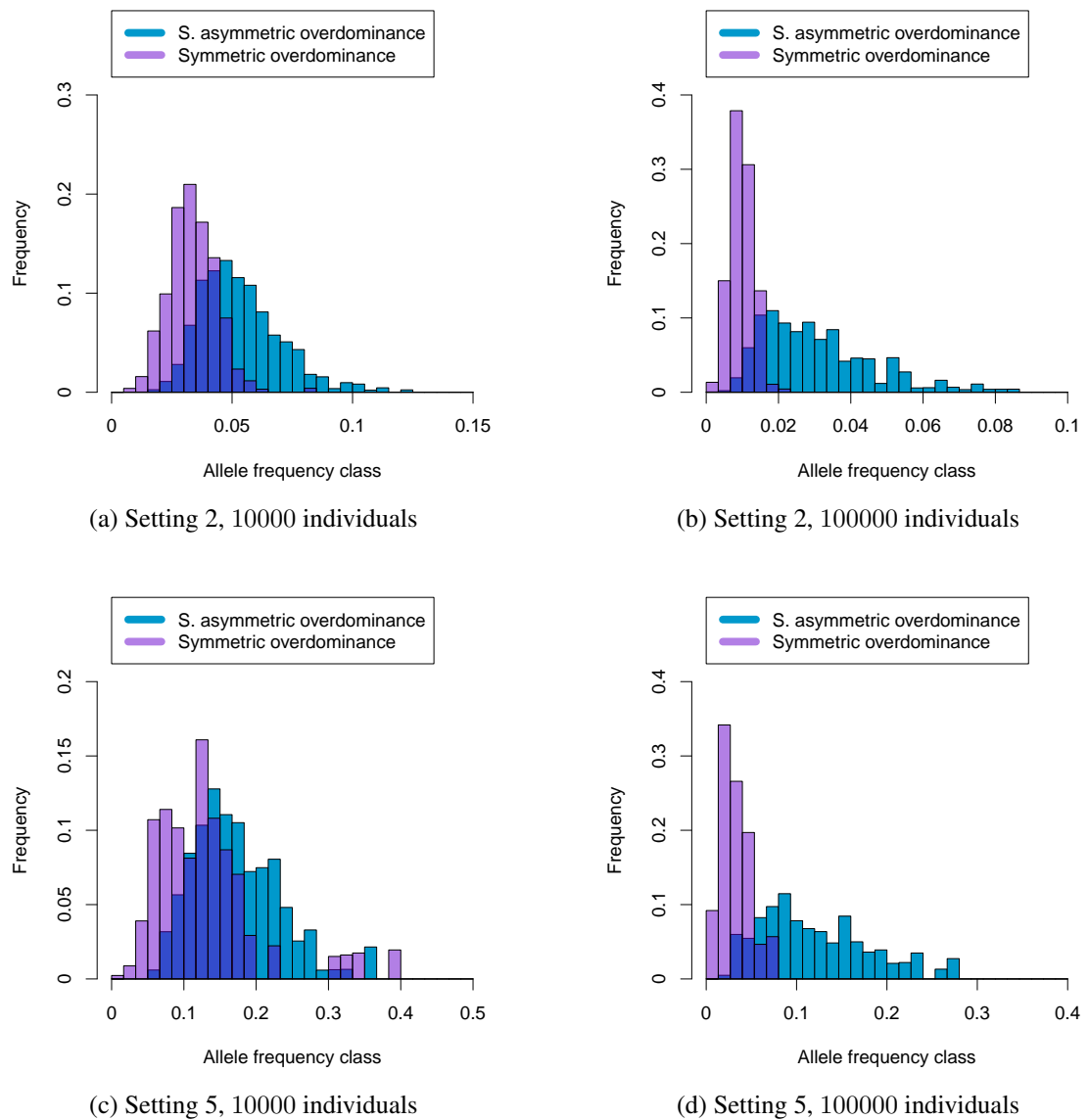


Figure B.6: Weighted frequency spectra for the SOD model and the SAsOD model for settings 2 and 5 and population sizes of 10000 and 100000.

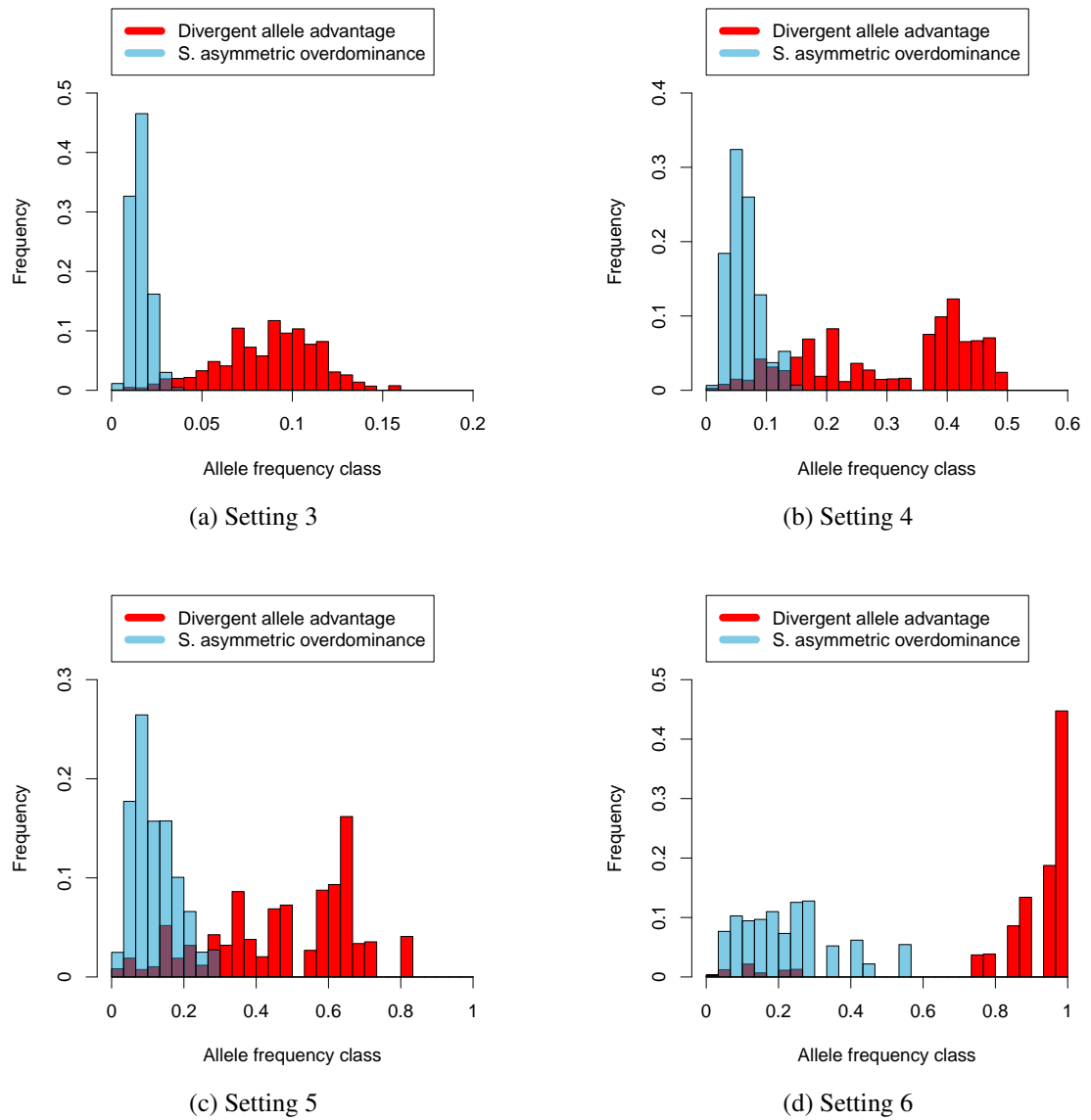


Figure B.7: Weighted frequency spectra for the DAA model and the SAsOD model for settings 3, 4, 5 and 6. The population size is always 100000.

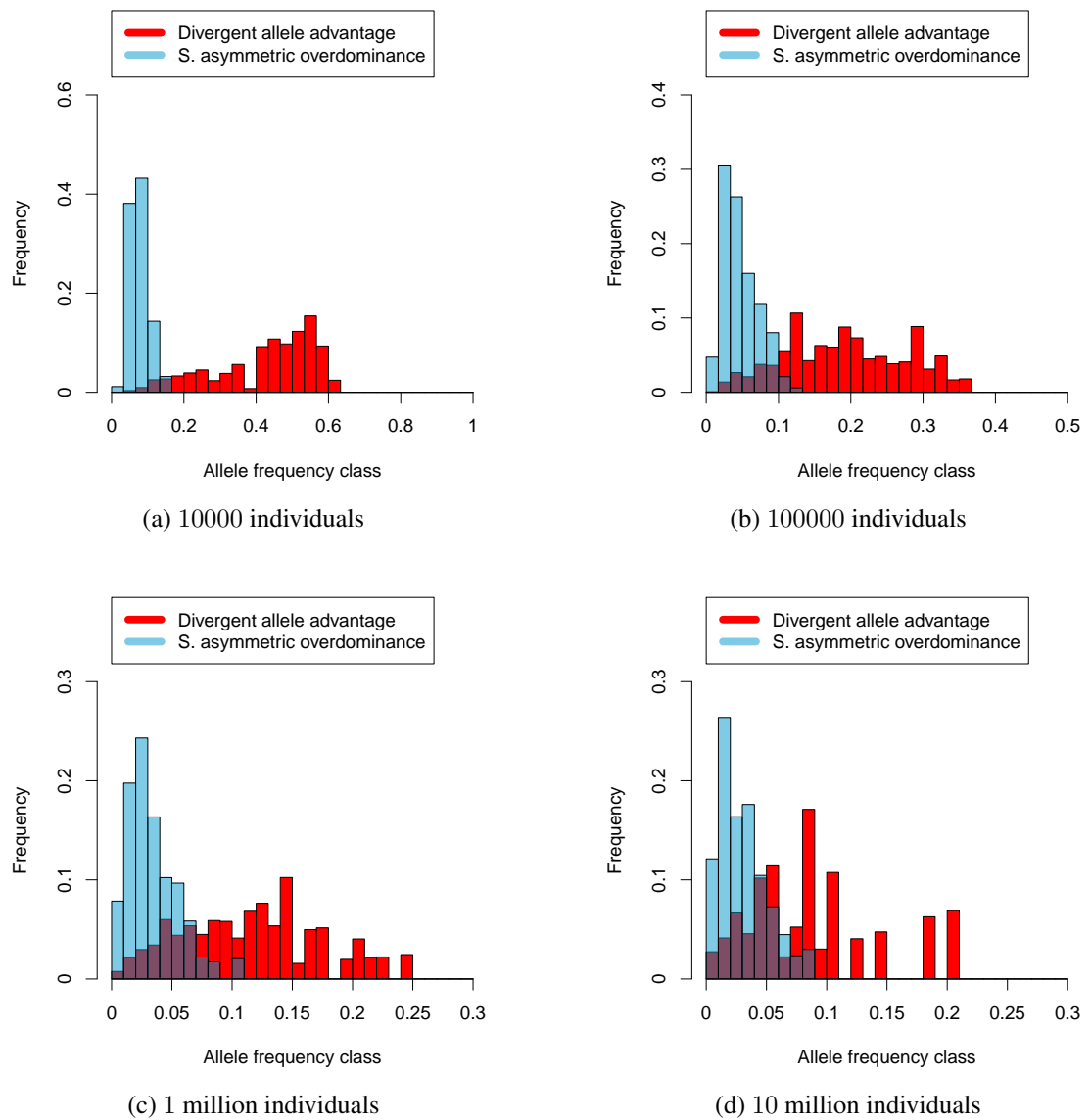


Figure B.8: Weighted frequency spectra for the DAA model and the SASOD model for setting 1 and different population sizes.

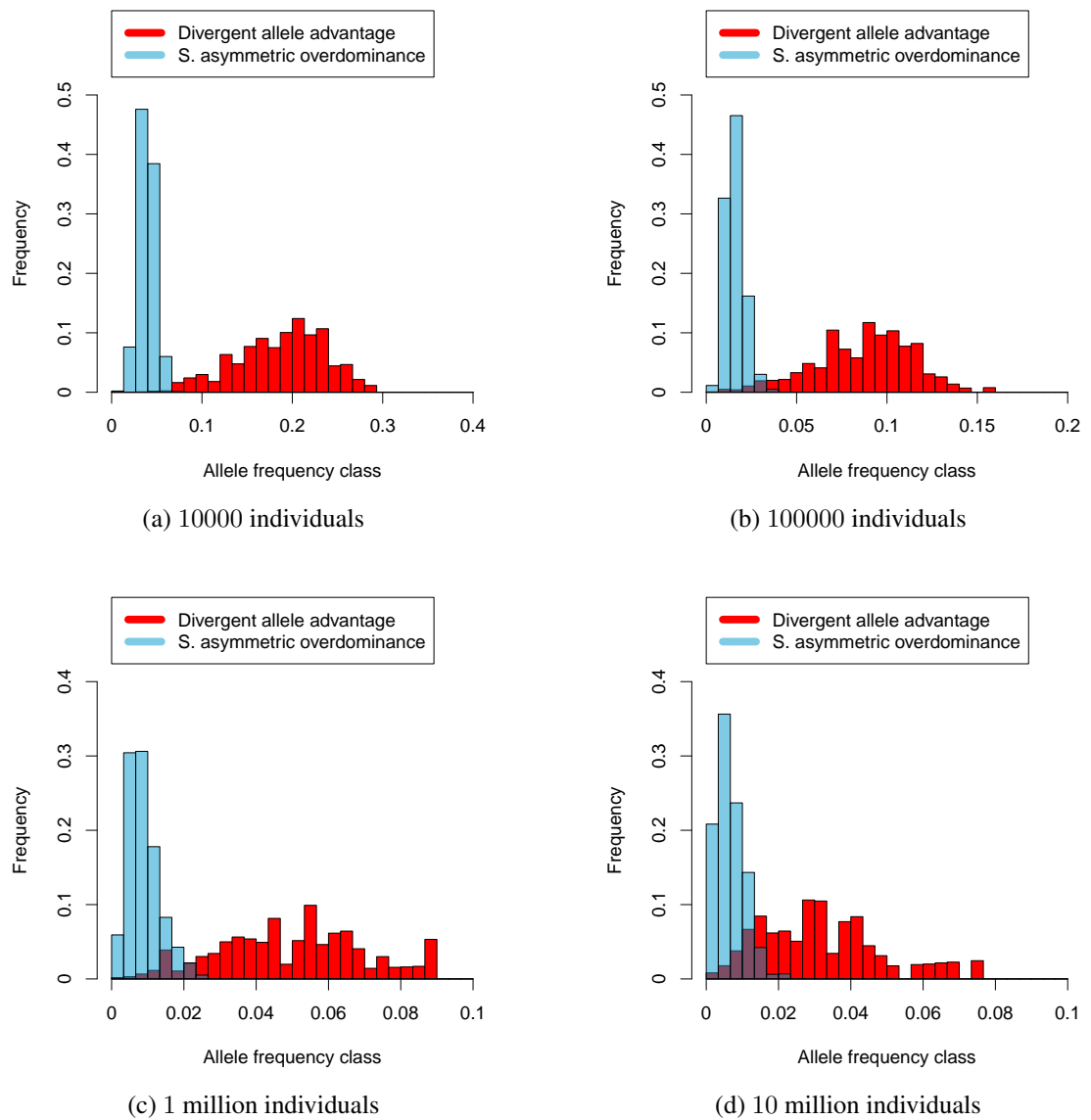


Figure B.9: Weighted frequency spectra for the DAA model and the SAsOD model for setting 3 and different population sizes.

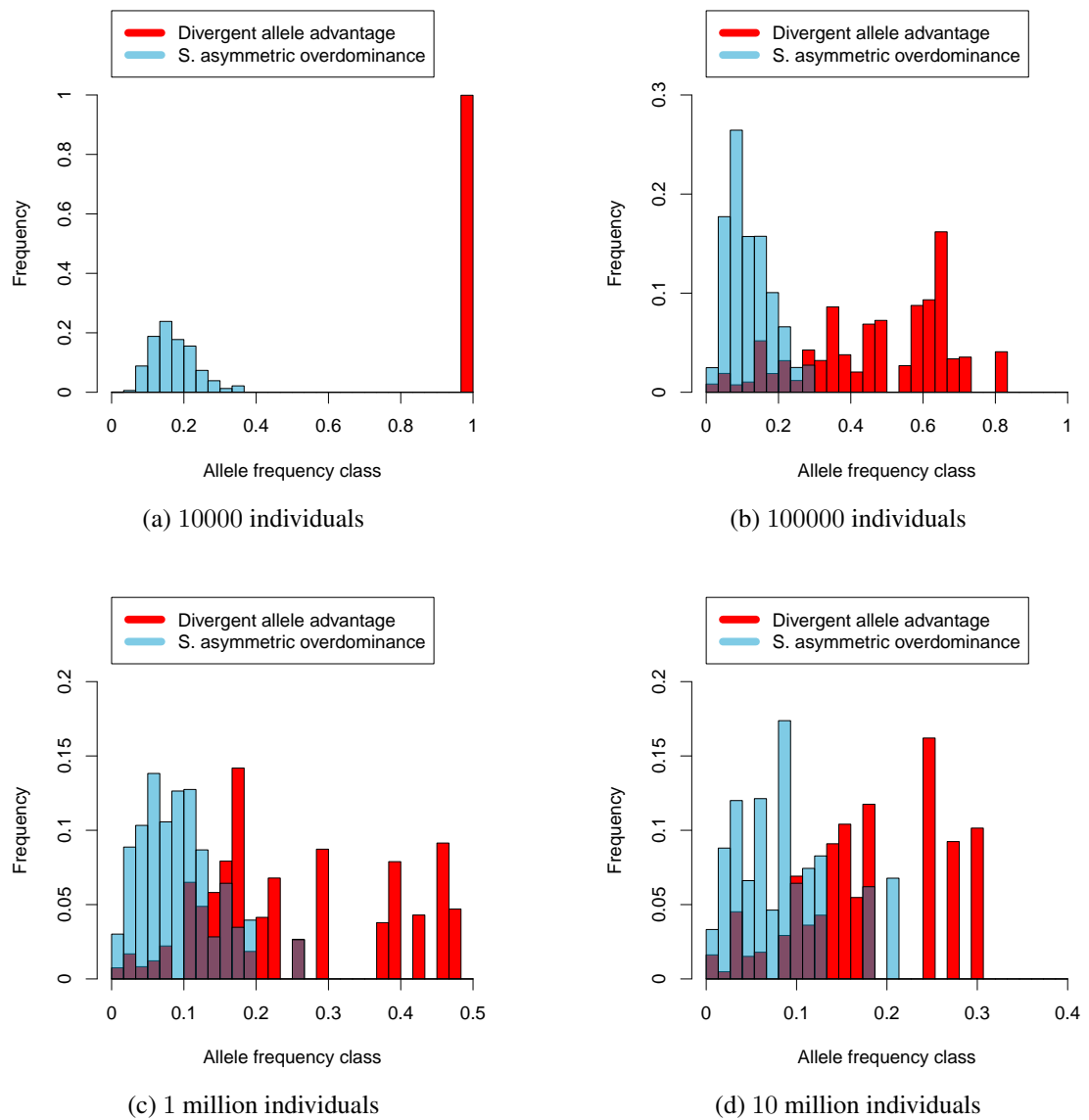


Figure B.10: Weighted frequency spectra for the DAA model and the SASOD model for setting 5 and different population sizes.

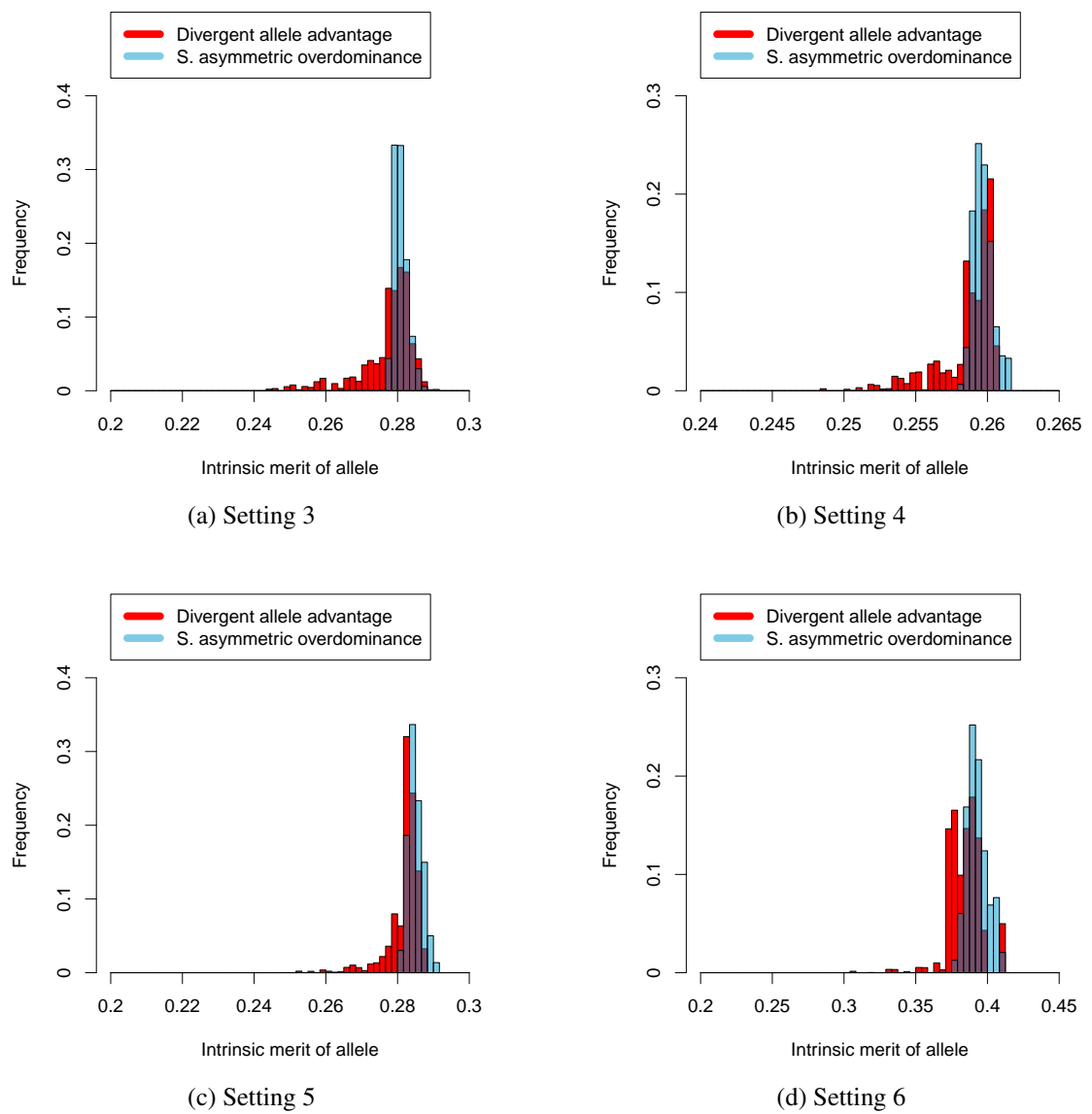


Figure B.11: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SASOD model) for setting 3 (top left), setting 4 (top right), setting 5 (bottom left) and setting 6 (bottom right). The population size was 100000.

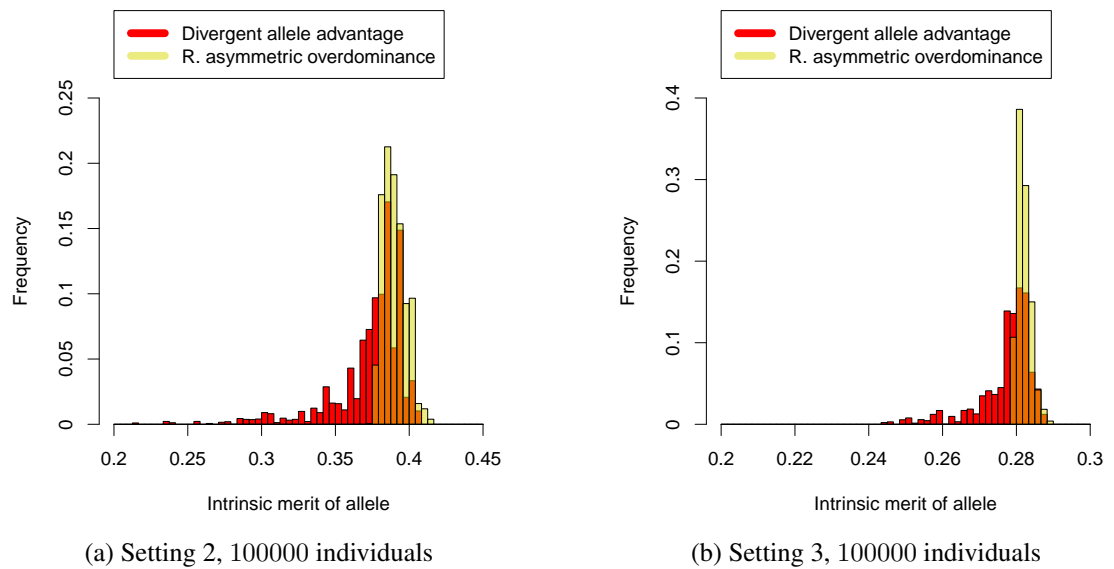


Figure B.12: Comparison between allele intrinsic merit distributions of the DAA and the RAsOD model.

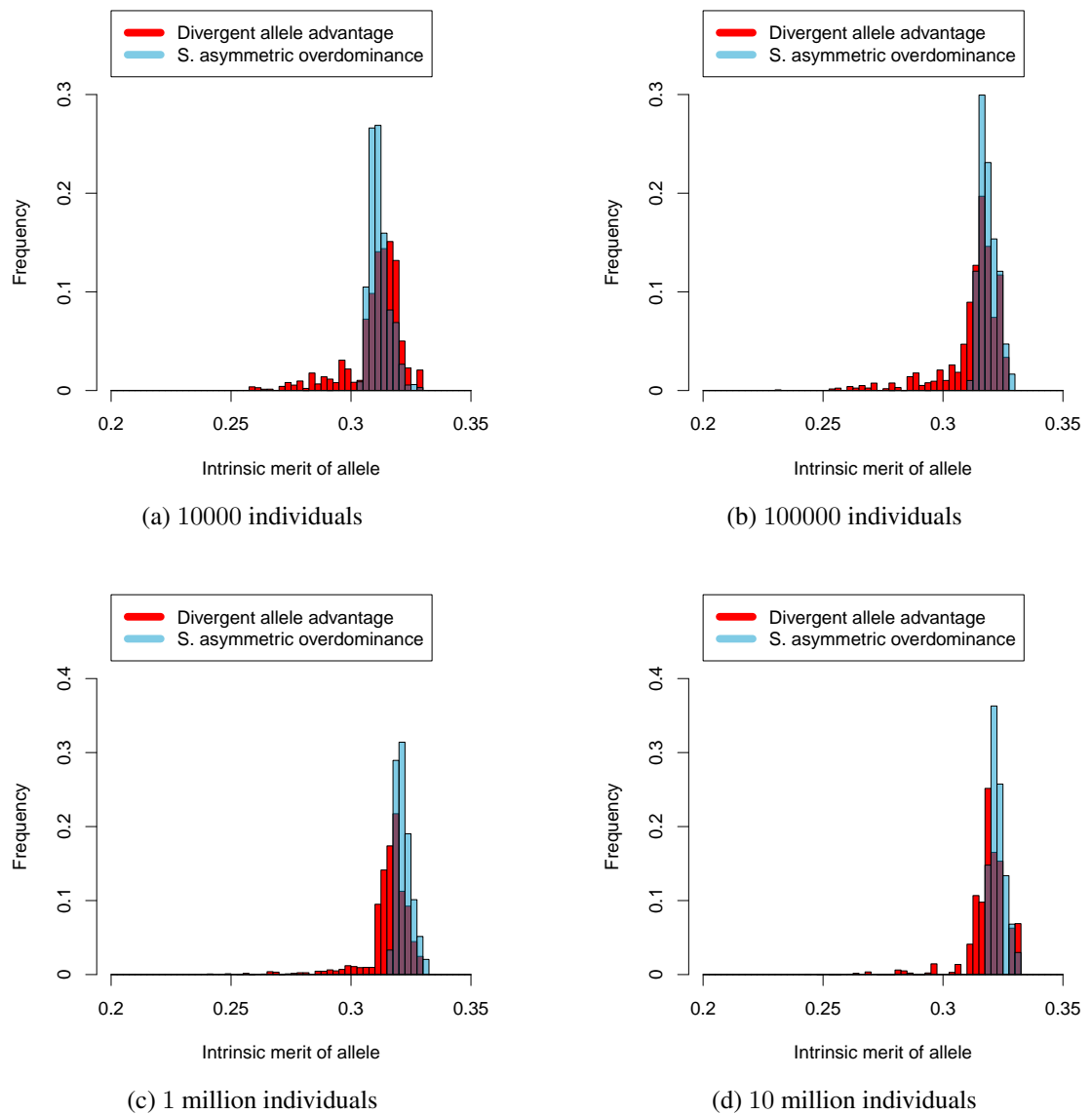


Figure B.13: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 1 and different population sizes.

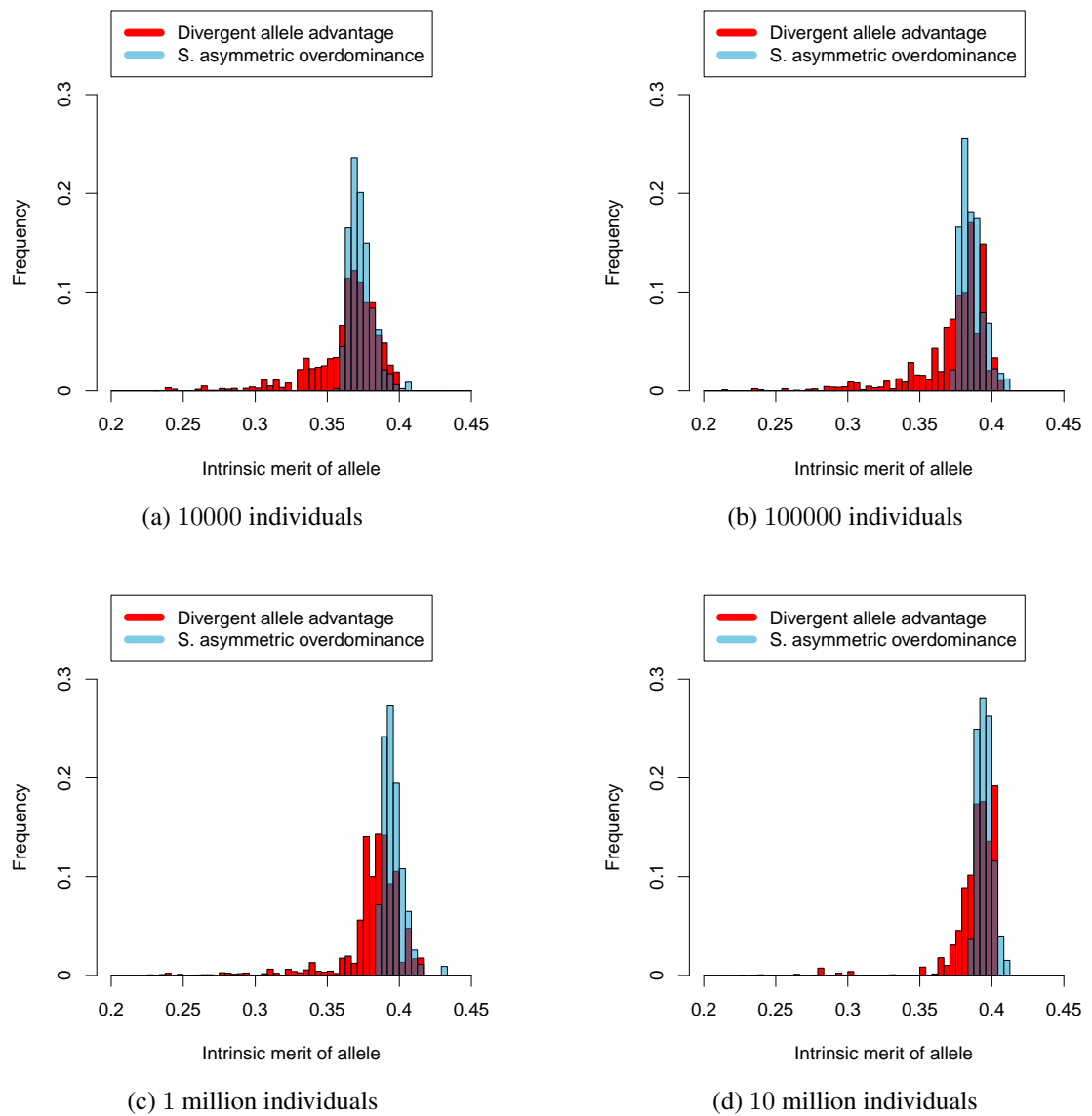


Figure B.14: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 2 and different population sizes.

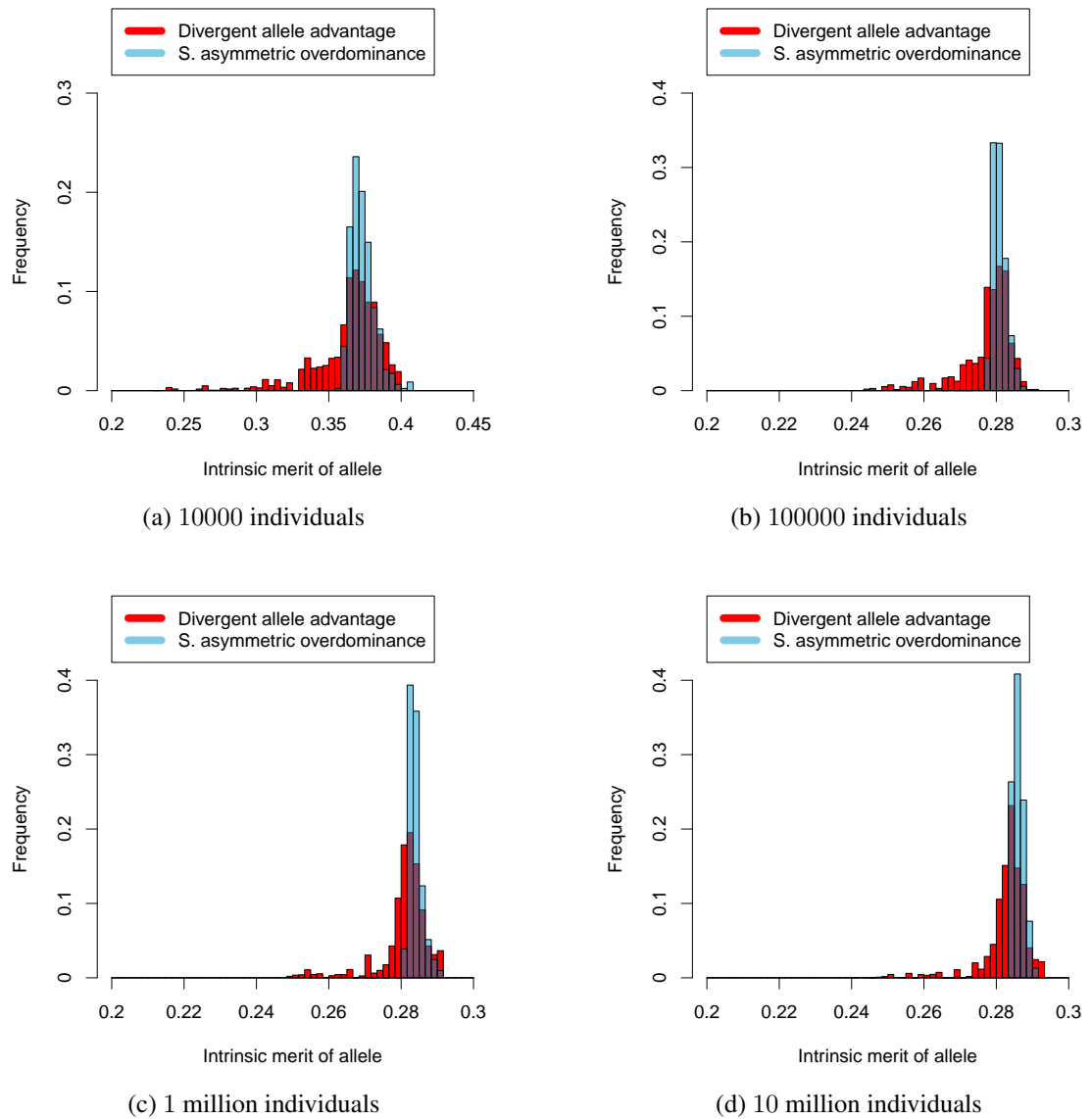


Figure B.15: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 3 and different population sizes.

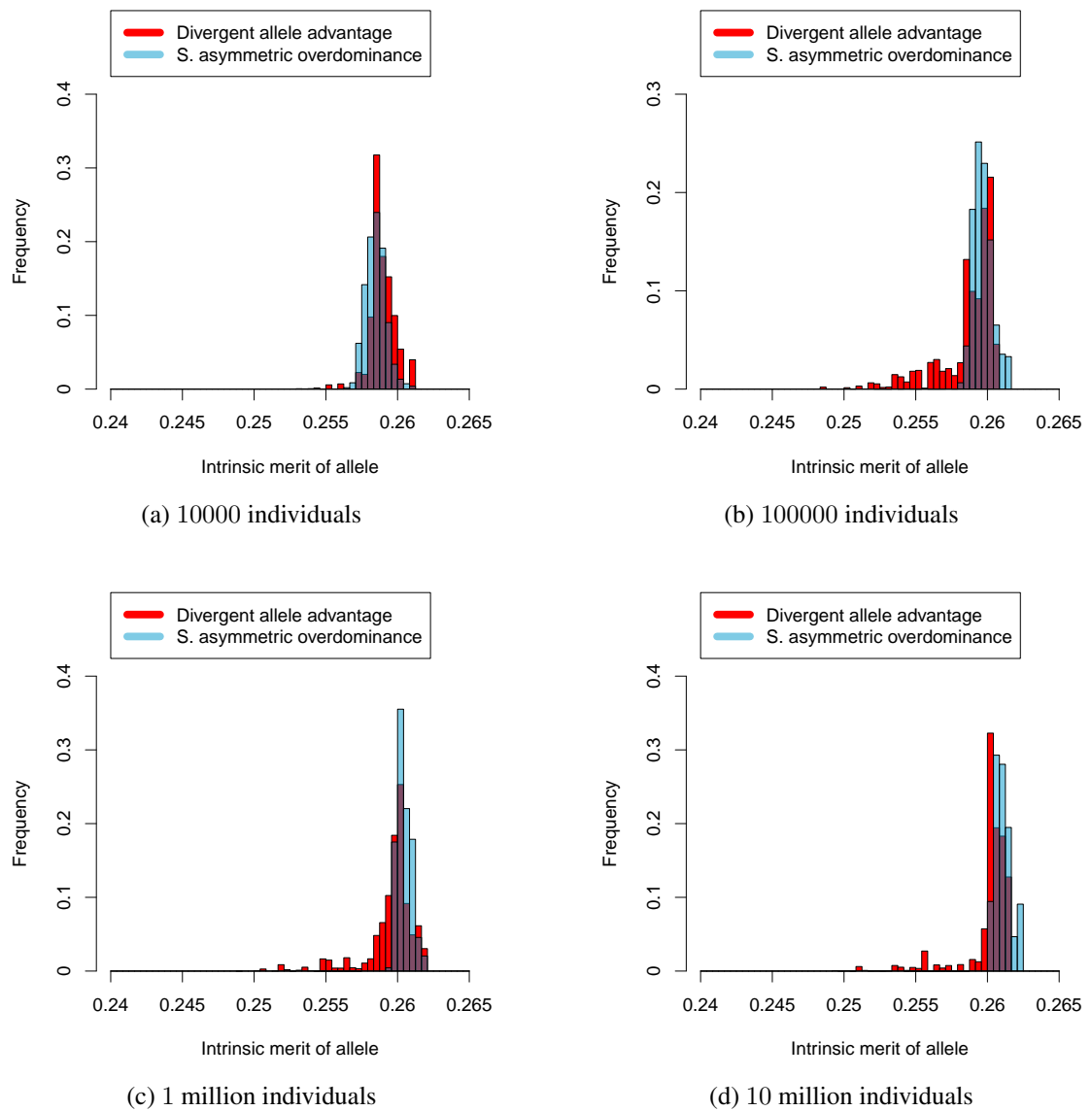


Figure B.16: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 4 and different population sizes.

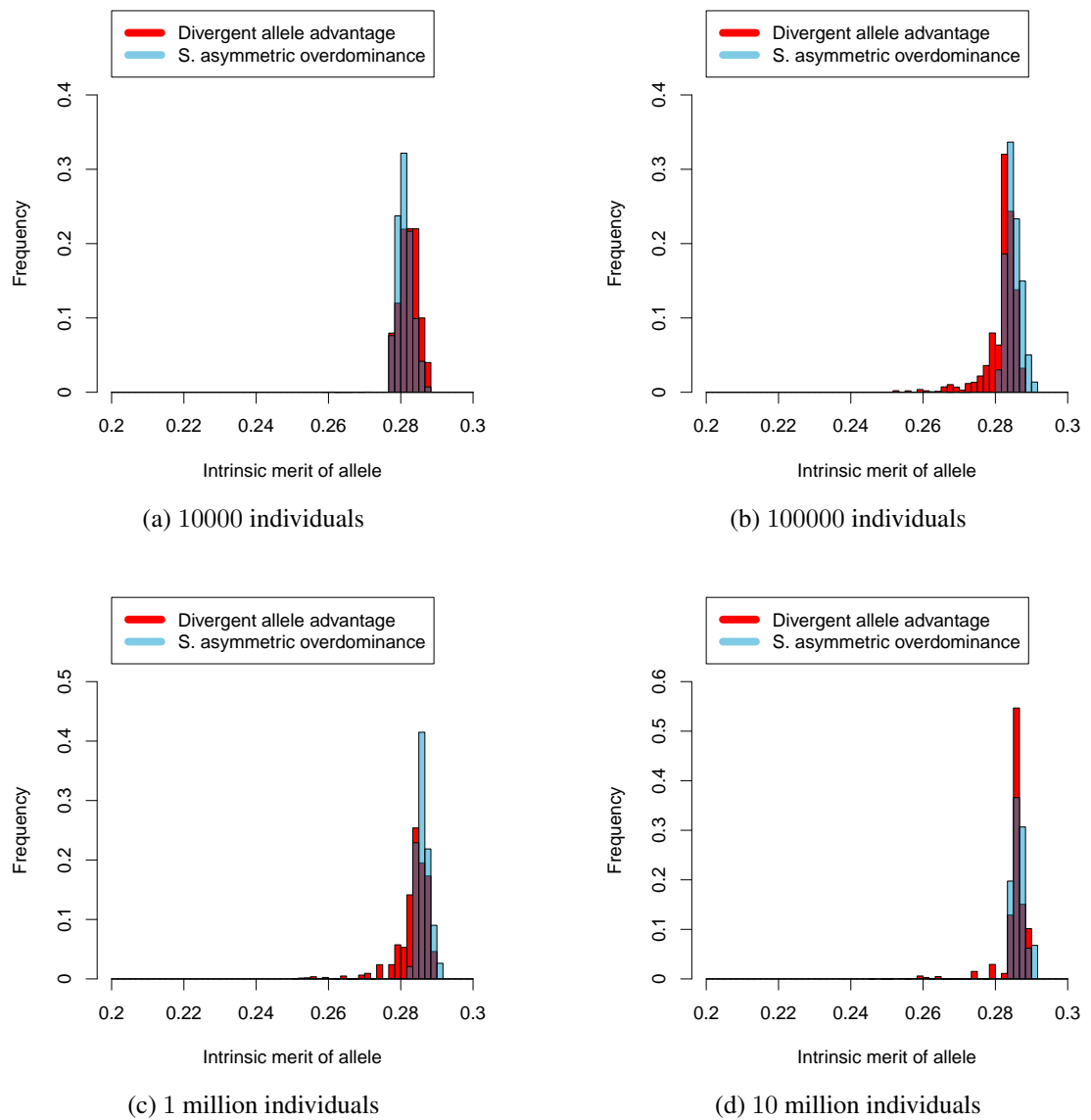


Figure B.17: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 5 and different population sizes.

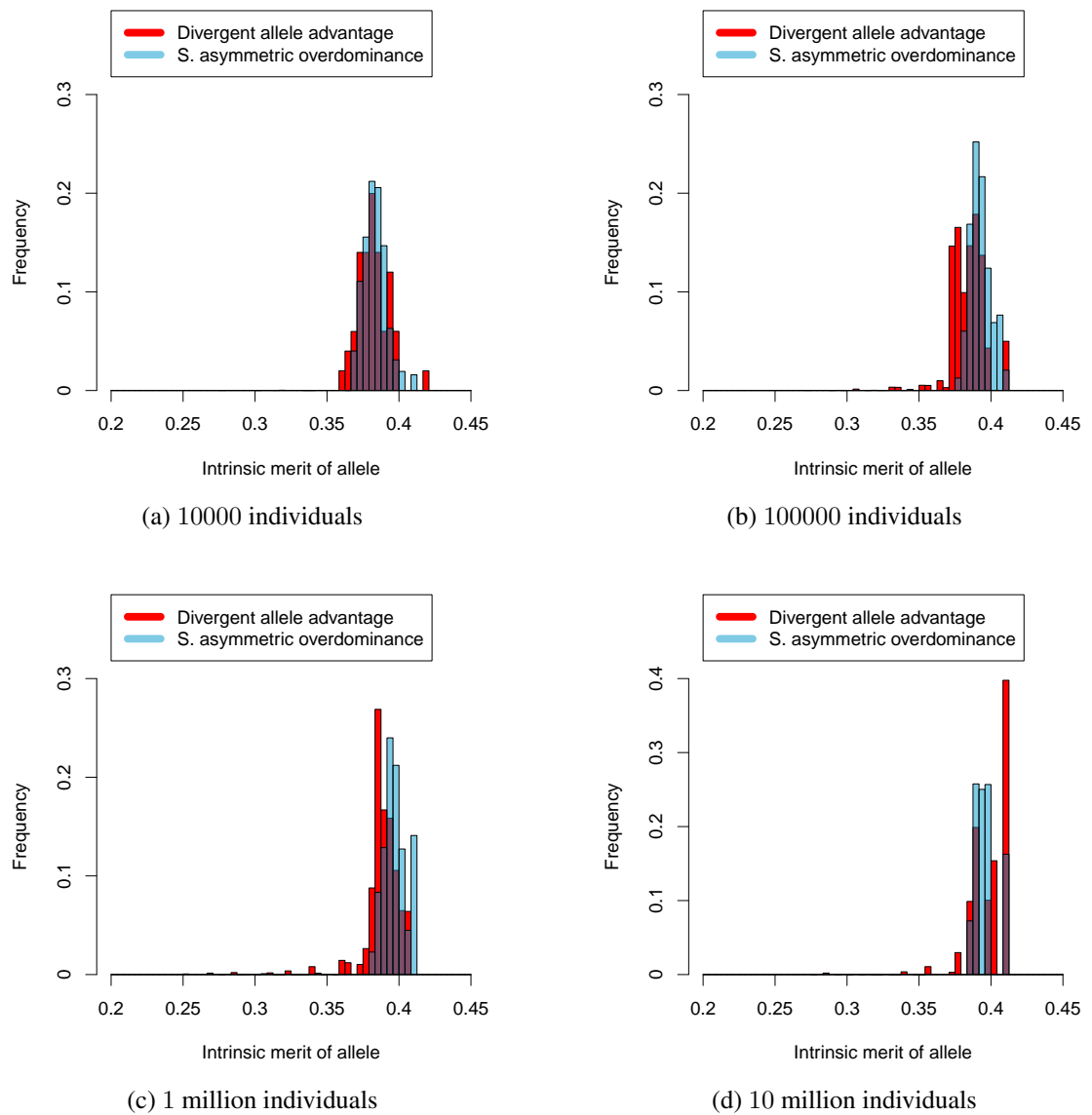
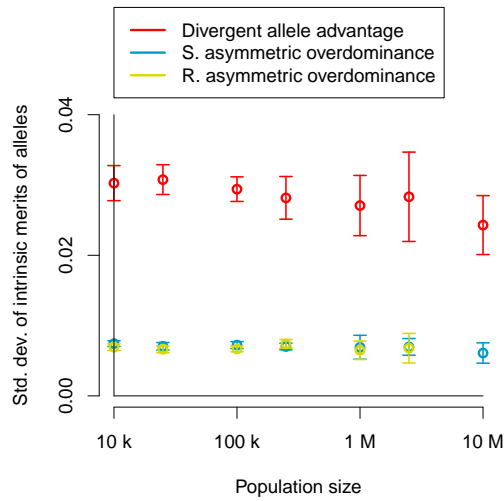
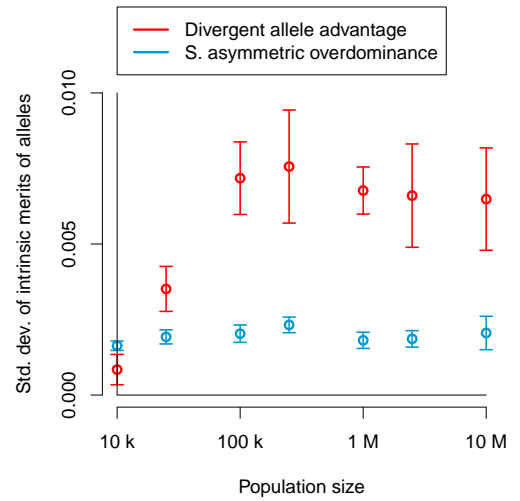


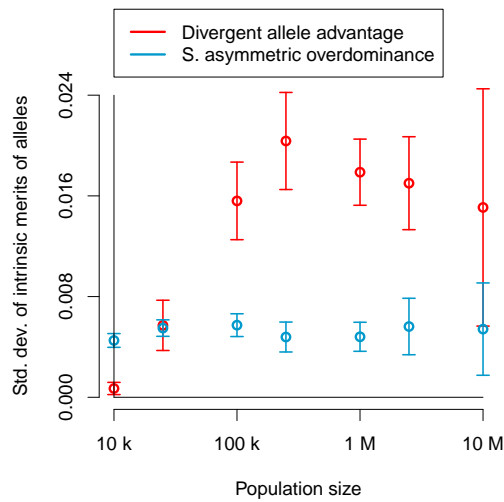
Figure B.18: Greater variation in intrinsic merit under DAA compared to a traditional asymmetric overdominance model (the SAsOD model) for setting 6 and different population sizes.



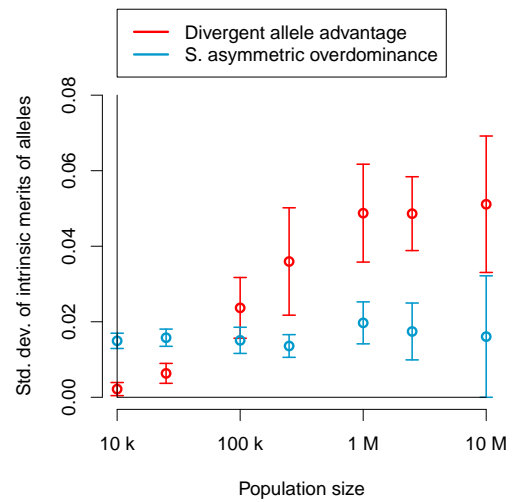
(a) Setting 3



(b) Setting 4



(c) Setting 5



(d) Setting 6

Figure B.19: Variation in allele intrinsic merit, measured as the standard deviation of the intrinsic merits of the final set of alleles. The alleles were weighted by their proportions. The DAA model is associated with approximately a 3x – 4x larger standard deviation than the asymmetric overdominance models for most settings and population sizes explored.

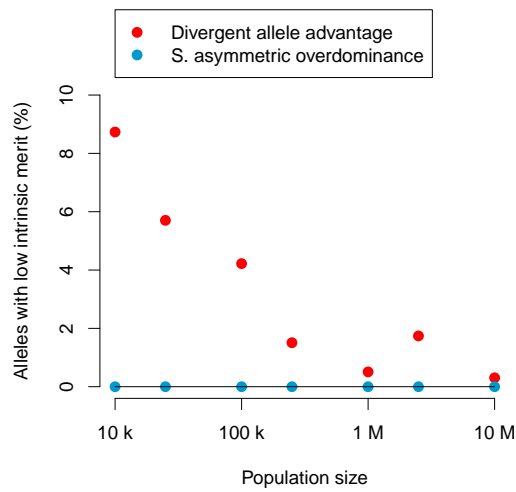
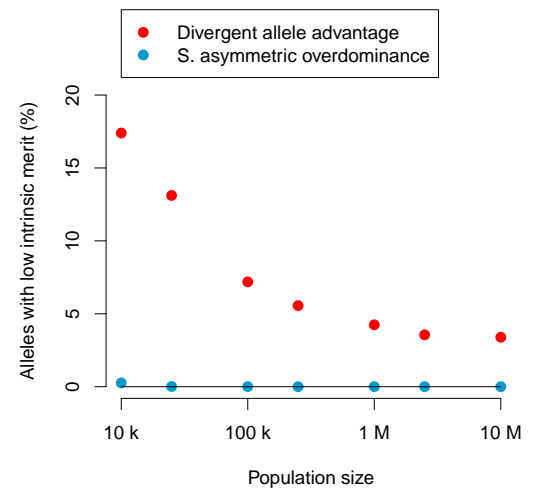
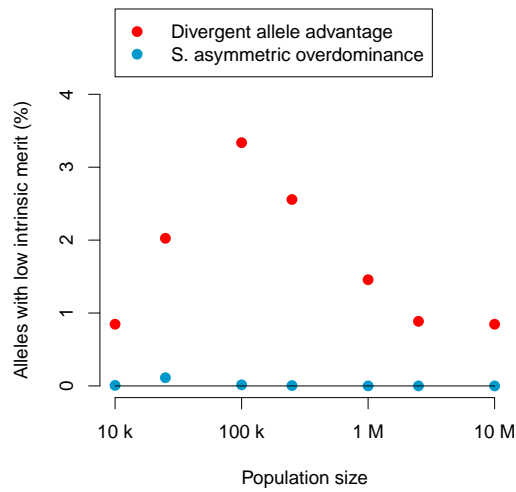
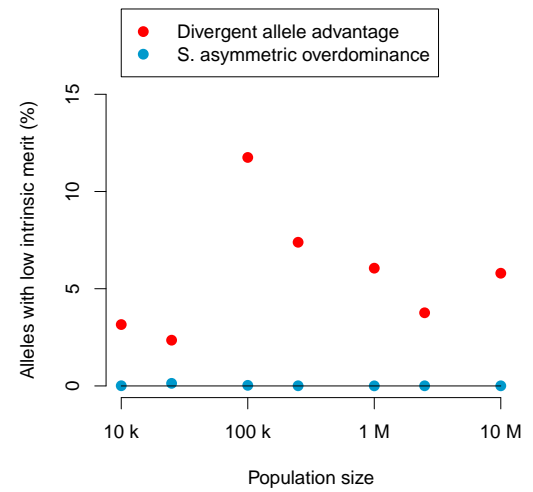
(a) Setting 3, $w_t = 0.25$ (b) Setting 3, $w_t = 0.265$ (c) Setting 4, $w_t = 0.253$ (d) Setting 4, $w_t = 0.256$

Figure B.20: Percentage of alleles below two different thresholds of allele intrinsic merit (w_t) for the DAA and SAsOD models and settings 3 and 4.

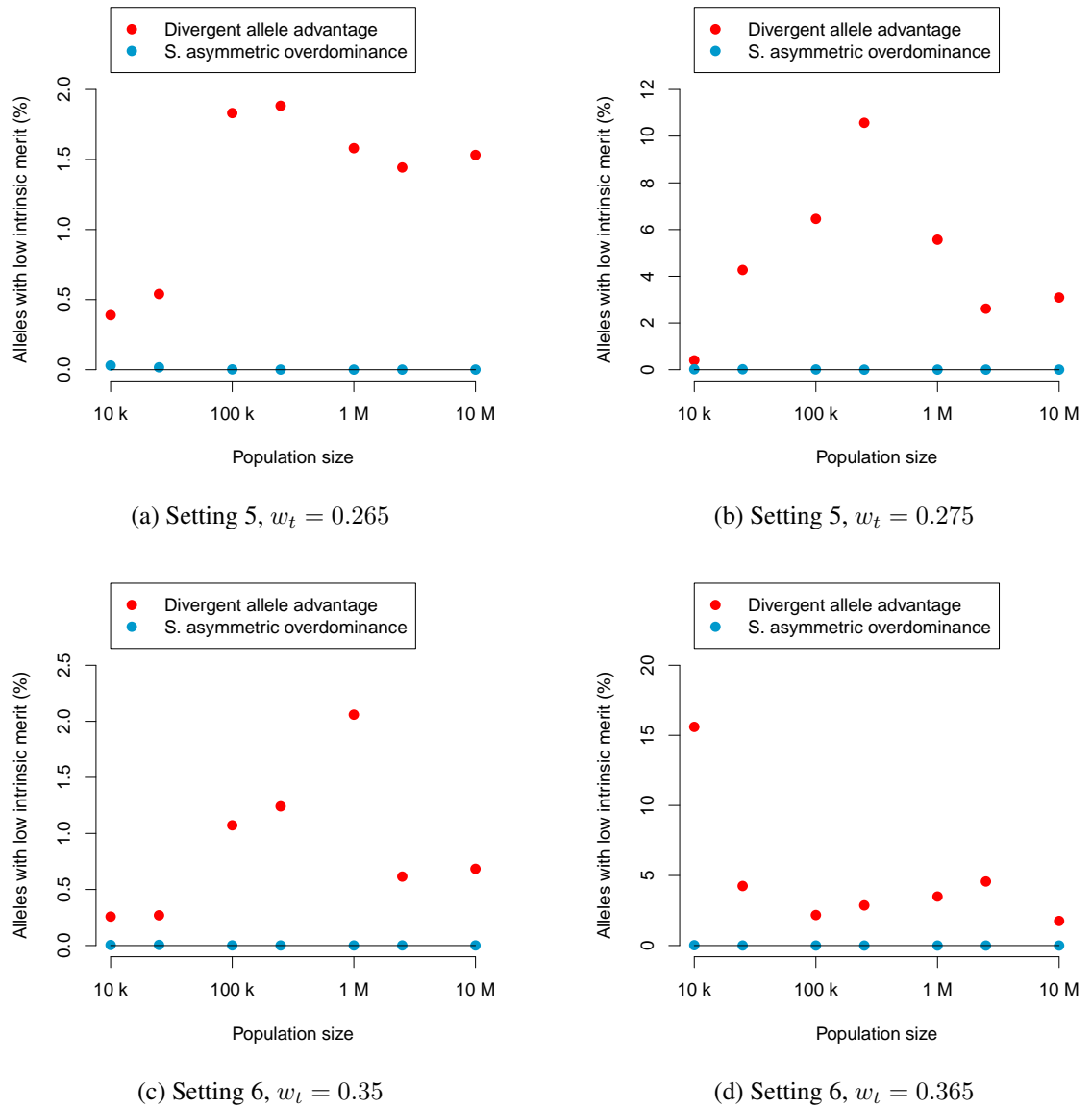


Figure B.21: Percentage of alleles below two different thresholds of allele intrinsic merit (w_t) for the DAA and SAsOD models and settings 5 and 6.

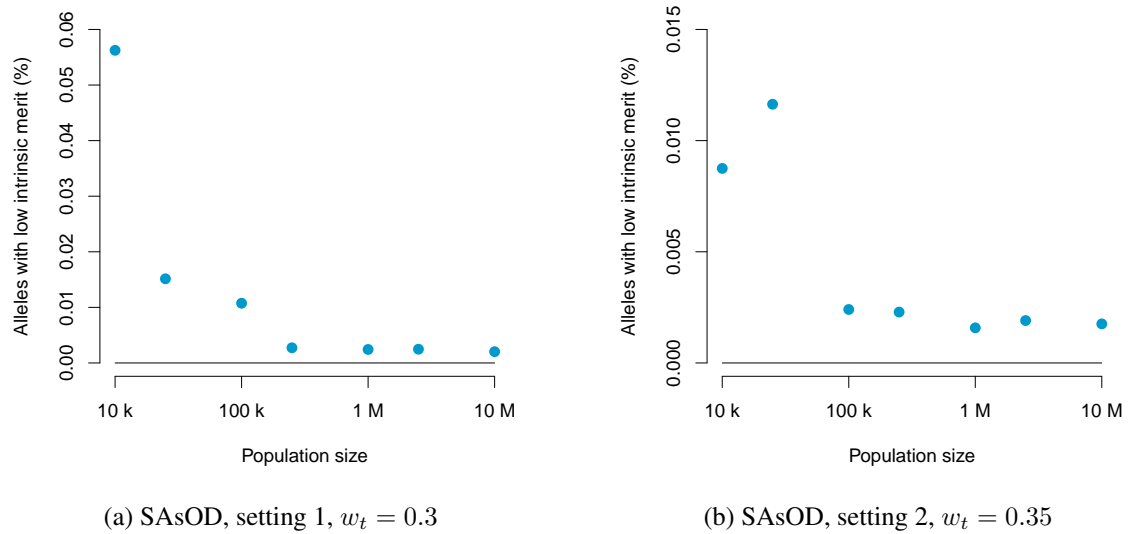


Figure B.22: Percentage of alleles below a thresholds of allele intrinsic merit of 0.3 (setting 1, left) and 0.35 (setting 2, right) for the SAsOD model.

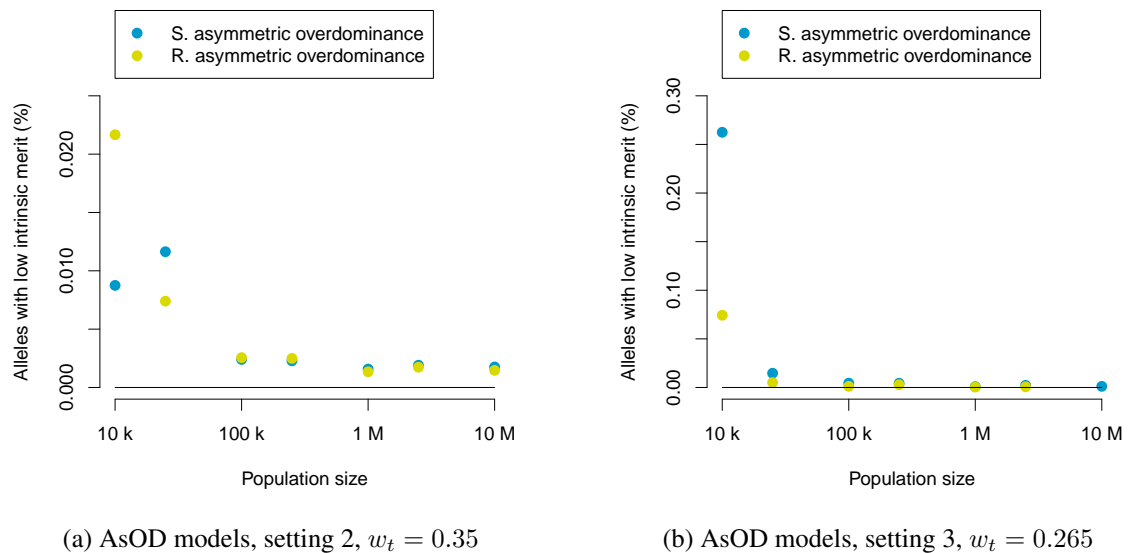


Figure B.23: Percentage of alleles below a thresholds of allele intrinsic merit of 0.35 (setting 2, left) and 0.265 (setting 3, right) for the SAsOD and the RAsOD model. For a population size of 10 million, only the SAsOD model was simulated.

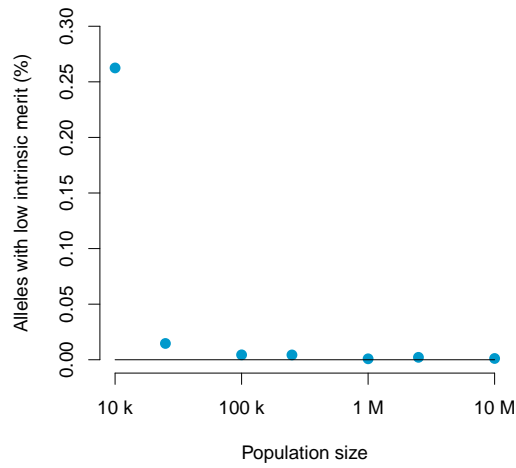
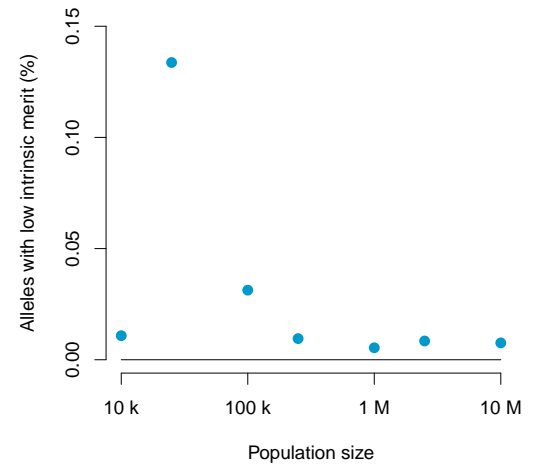
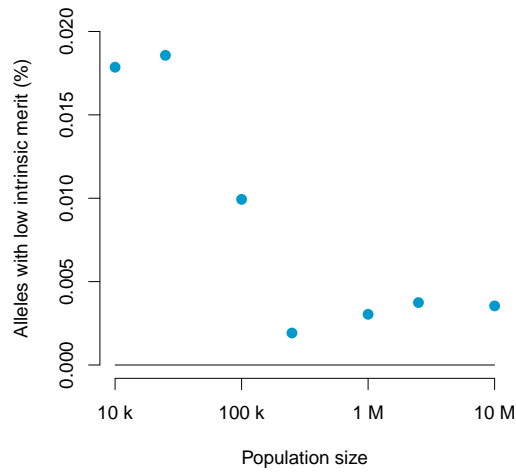
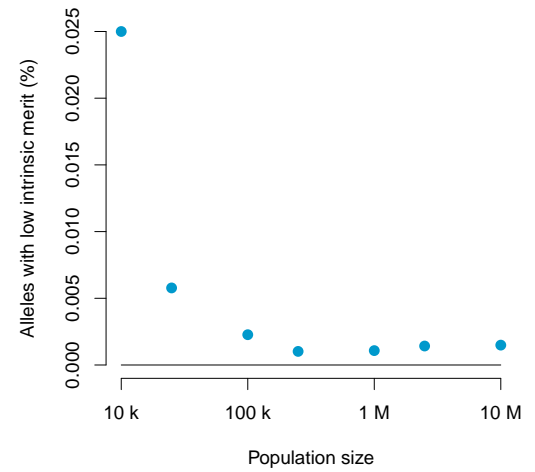
(a) SAsOD, setting 3, $w_t = 0.265$ (b) SAsOD, setting 4, $w_t = 0.256$ (c) SAsOD, setting 5, $w_t = 0.275$ (d) SAsOD, setting 6, $w_t = 0.365$

Figure B.24: Percentage of alleles below a thresholds of allele intrinsic merit of 0.265 (setting 3, top left), 0.256 (setting 4, top right), 0.275 (setting 5, bottom left) and 0.365 (setting 6, bottom right) for the SAsOD model.

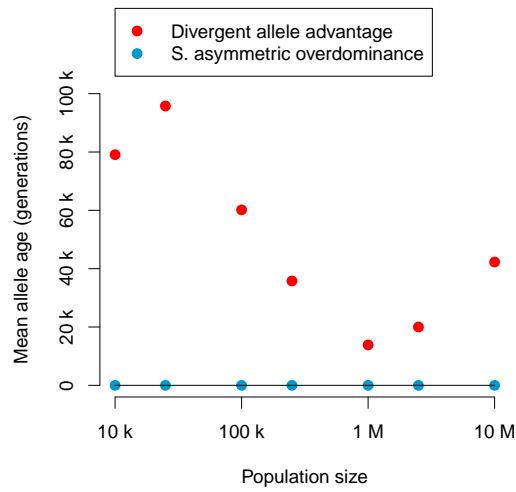
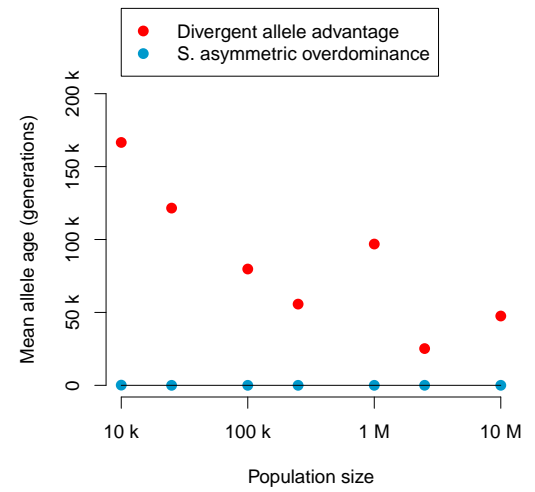
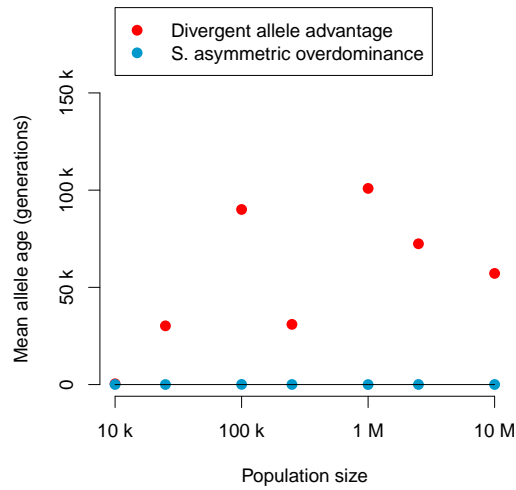
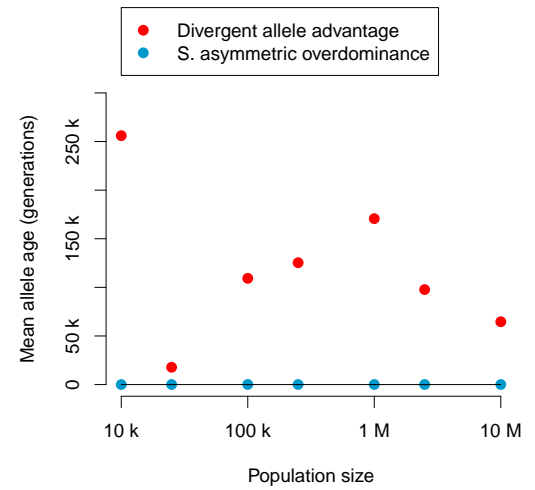
(a) Setting 3, $w_t = 0.25$ (b) Setting 3, $w_t = 0.265$ (c) Setting 4, $w_t = 0.253$ (d) Setting 4, $w_t = 0.256$

Figure B.25: Mean age of alleles below two different thresholds of allele intrinsic merit (w_t) for the DAA and SAsOD models and settings 3 and 4.

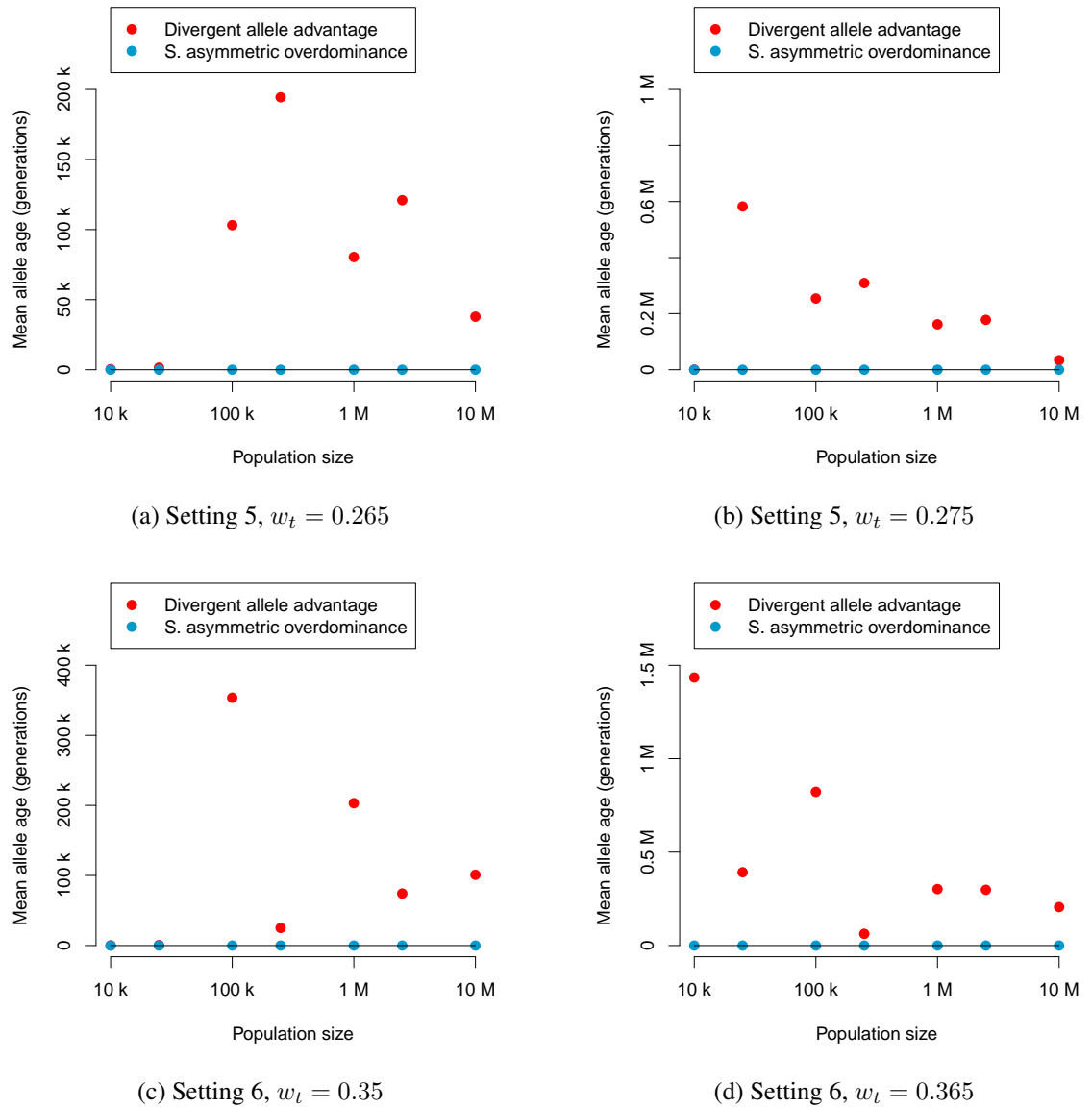


Figure B.26: Mean age of alleles below two different thresholds of allele intrinsic merit (w_t) for the DAA and SAsOD models and settings 5 and 6.

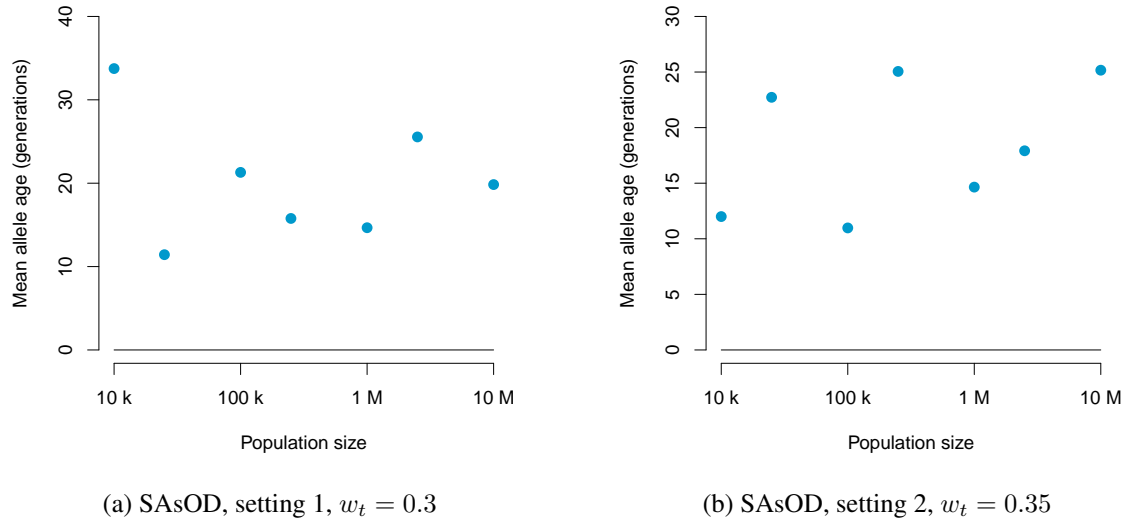


Figure B.27: Mean age of alleles below a thresholds of allele intrinsic merit of 0.3 (setting 1, left) and 0.35 (setting 2, right) for the SAsOD model.

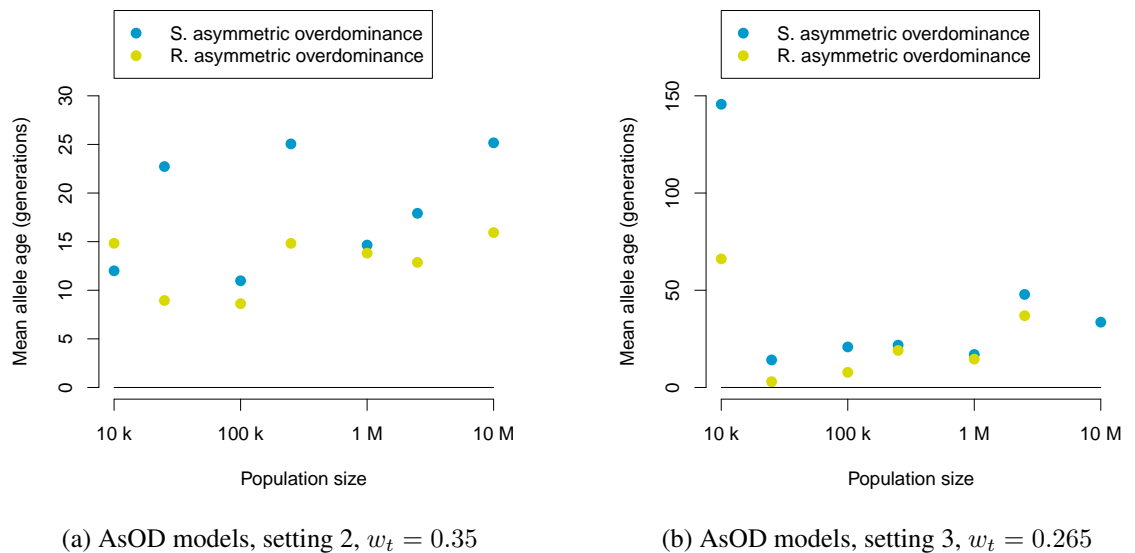


Figure B.28: Mean age of alleles below a thresholds of allele intrinsic merit of 0.35 (setting 2, left) and 0.265 (setting 3, right) for the SAsOD and the RAsOD model. For a population size of 10 million, only the SAsOD model was simulated.

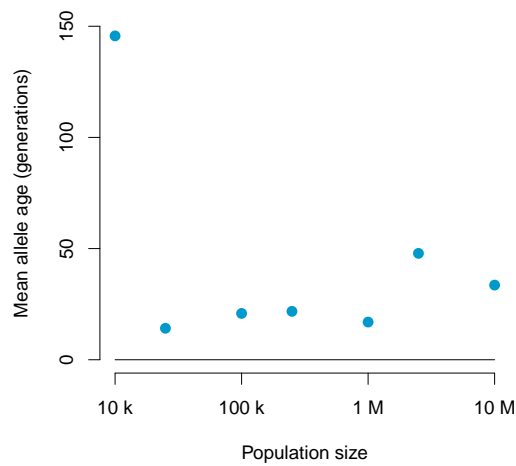
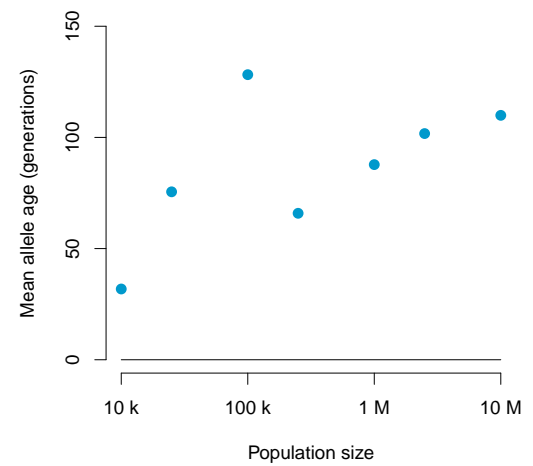
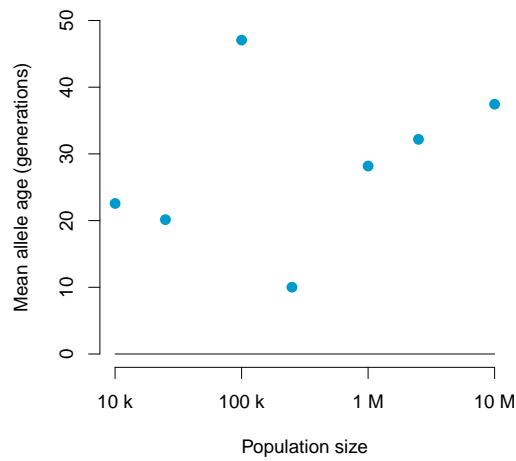
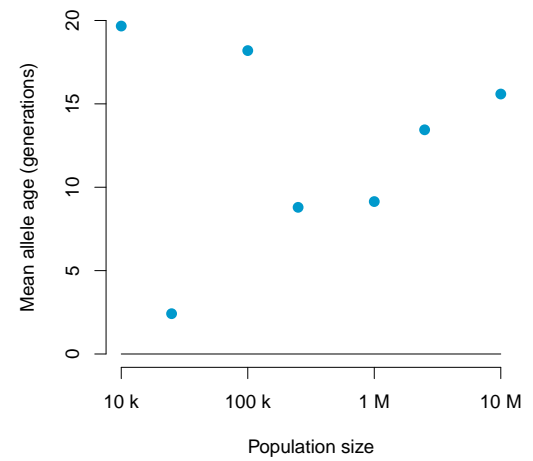
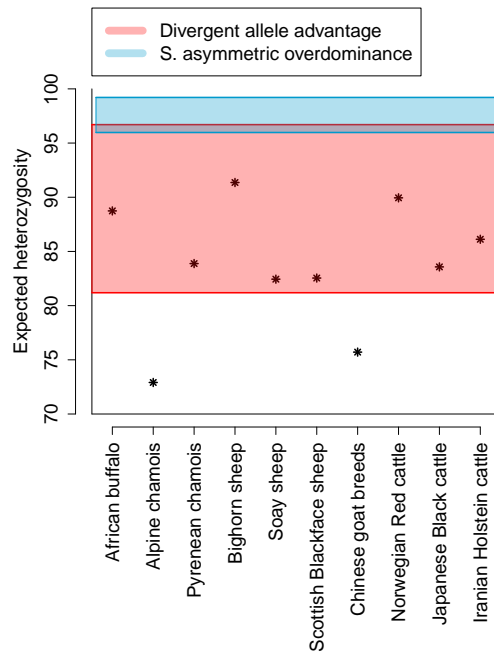
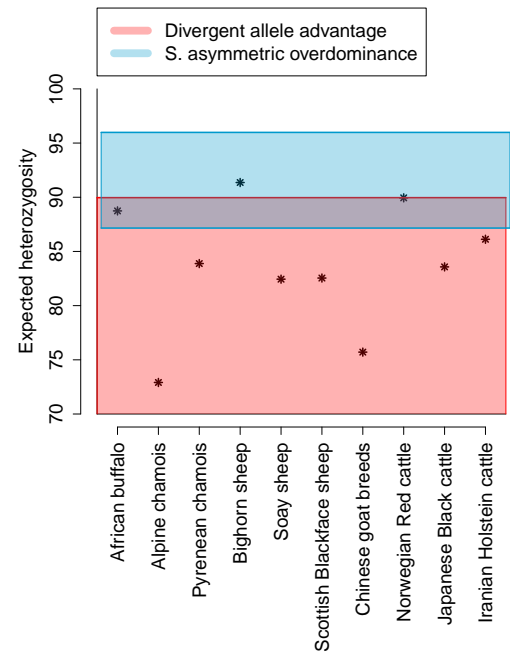
(a) SAsOD, setting 3, $w_t = 0.265$ (b) SAsOD, setting 4, $w_t = 0.256$ (c) SAsOD, setting 5, $w_t = 0.275$ (d) SAsOD, setting 6, $w_t = 0.365$

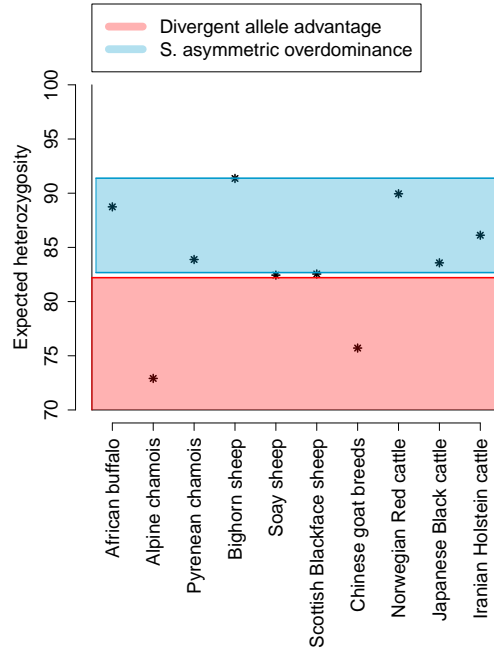
Figure B.29: Mean age of alleles below a thresholds of allele intrinsic merit of 0.265 (setting 3, top left), 0.256 (setting 4, top right), 0.275 (setting 5, bottom left) and 0.365 (setting 6, bottom right) for the SAsOD model.



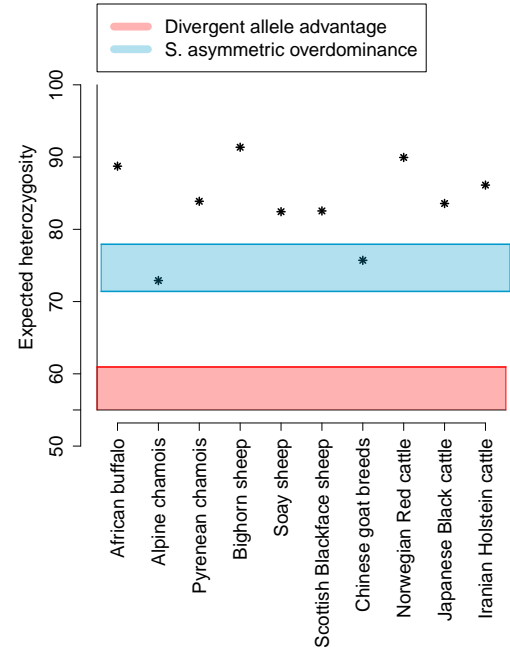
(a) Setting 3



(b) Setting 4

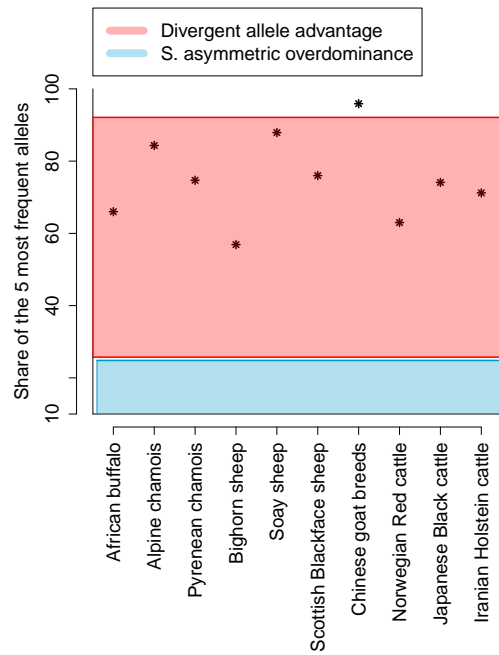


(c) Setting 5

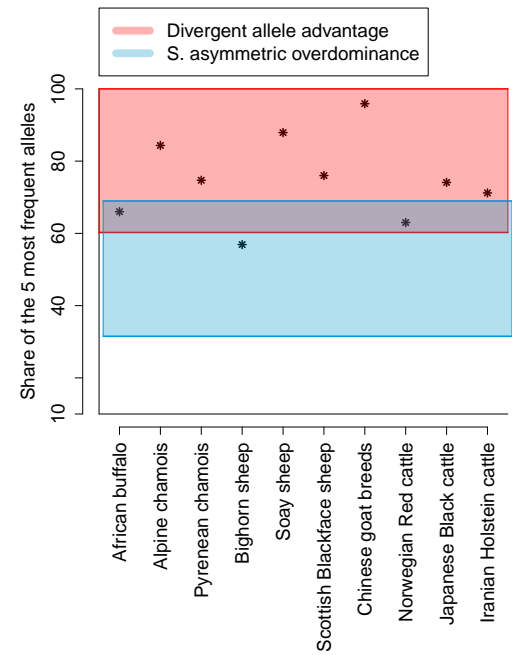


(d) Setting 6

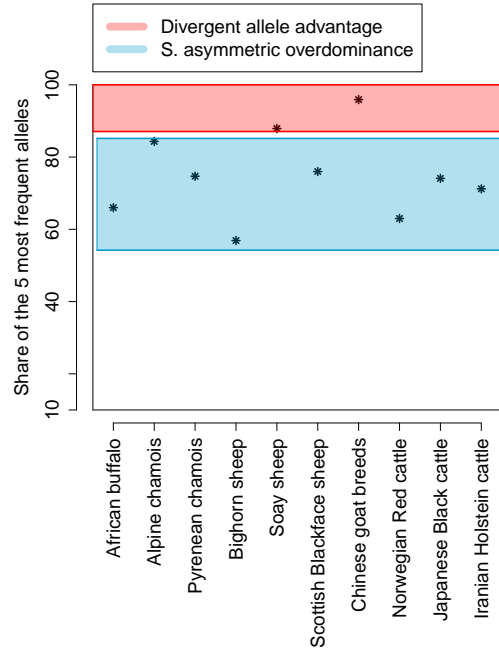
Figure B.30: Comparison to expected heterozygosities calculated from published allele frequencies for the DAA and SAsOD models and settings 3, 4, 5 and 6.



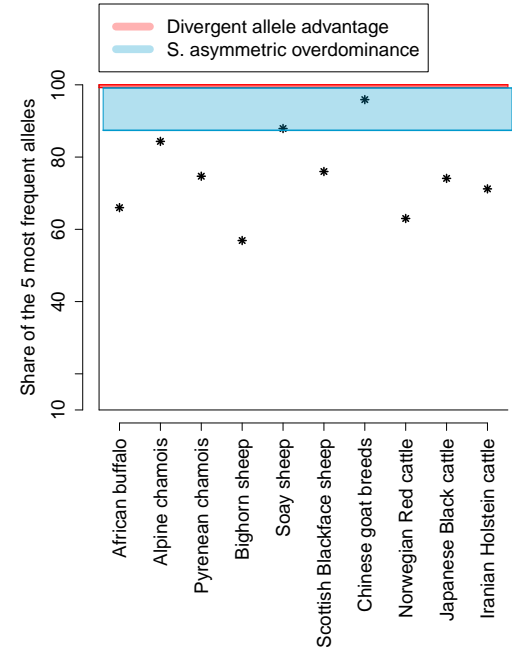
(a) Setting 3



(b) Setting 4

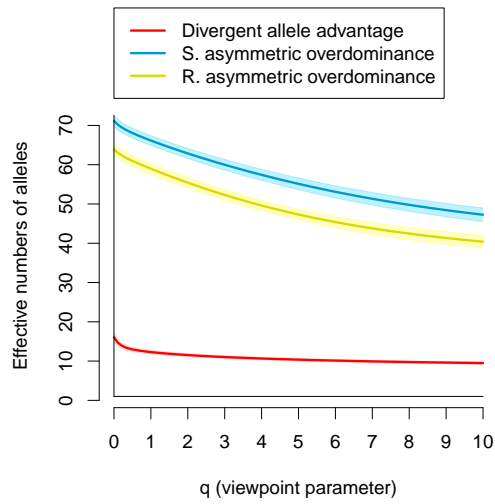


(c) Setting 5

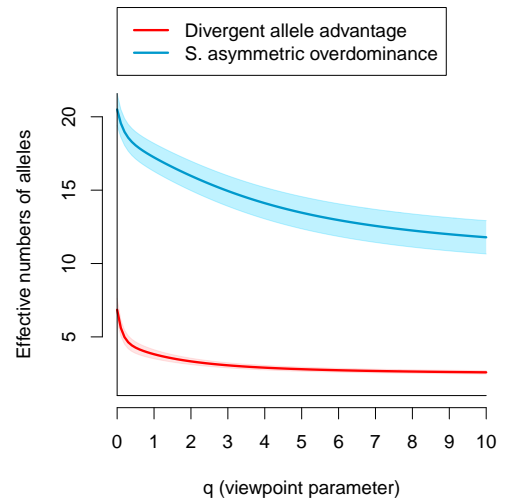


(d) Setting 6

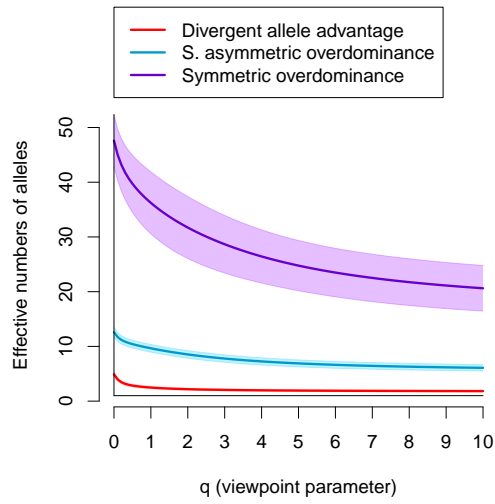
Figure B.31: Comparison to the share of the 5 most frequent alleles calculated from published allele frequencies for the DAA and SAsOD models and settings 3, 4, 5 and 6.



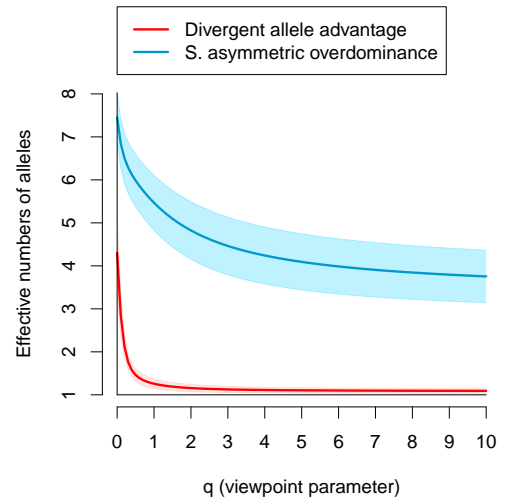
(a) Setting 3



(b) Setting 4



(c) Setting 5



(d) Setting 6

Figure B.32: Diversity profiles (naive diversity) for different overdominance models and settings 3, 4, 5 and 6. The population size was 100000.

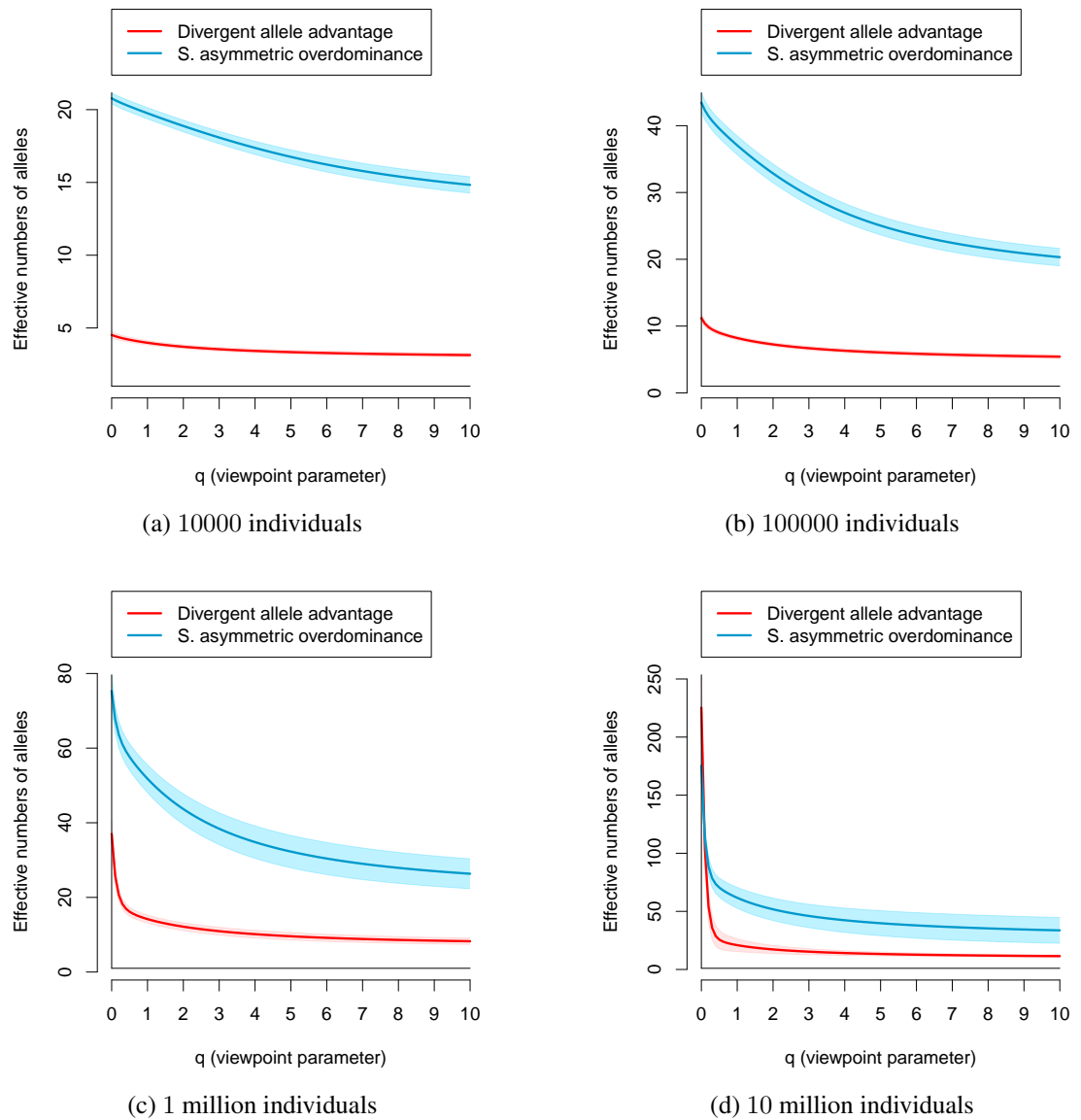


Figure B.33: Diversity profiles (naive diversity) for the DAA and SAsOD models, setting 2 and population sizes of 10000, 100000, 1 million and 10 million.

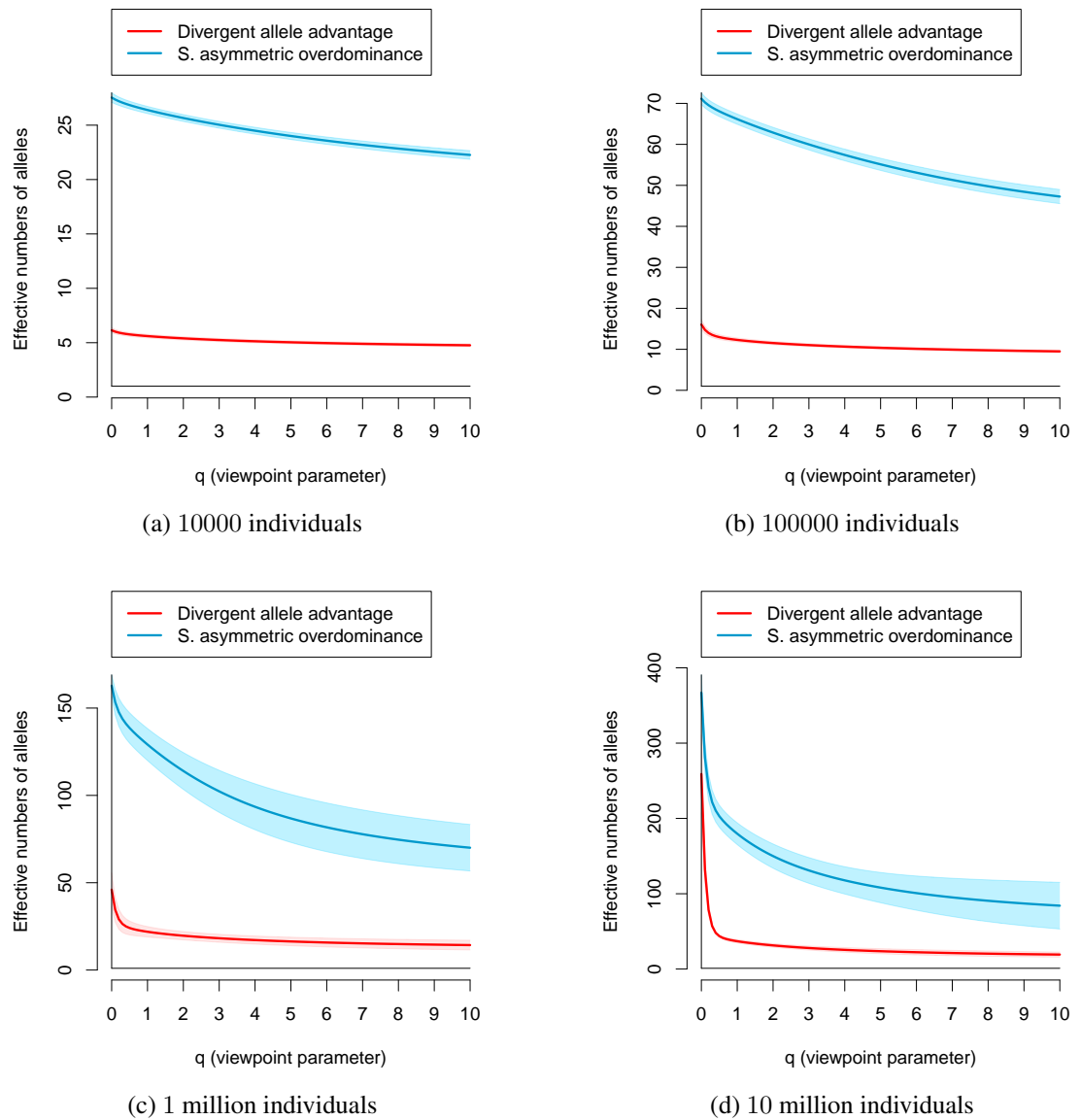


Figure B.34: Diversity profiles (naive diversity) for the DAA and SAsOD models, setting 3 and population sizes of 10000, 100000, 1 million and 10 million.

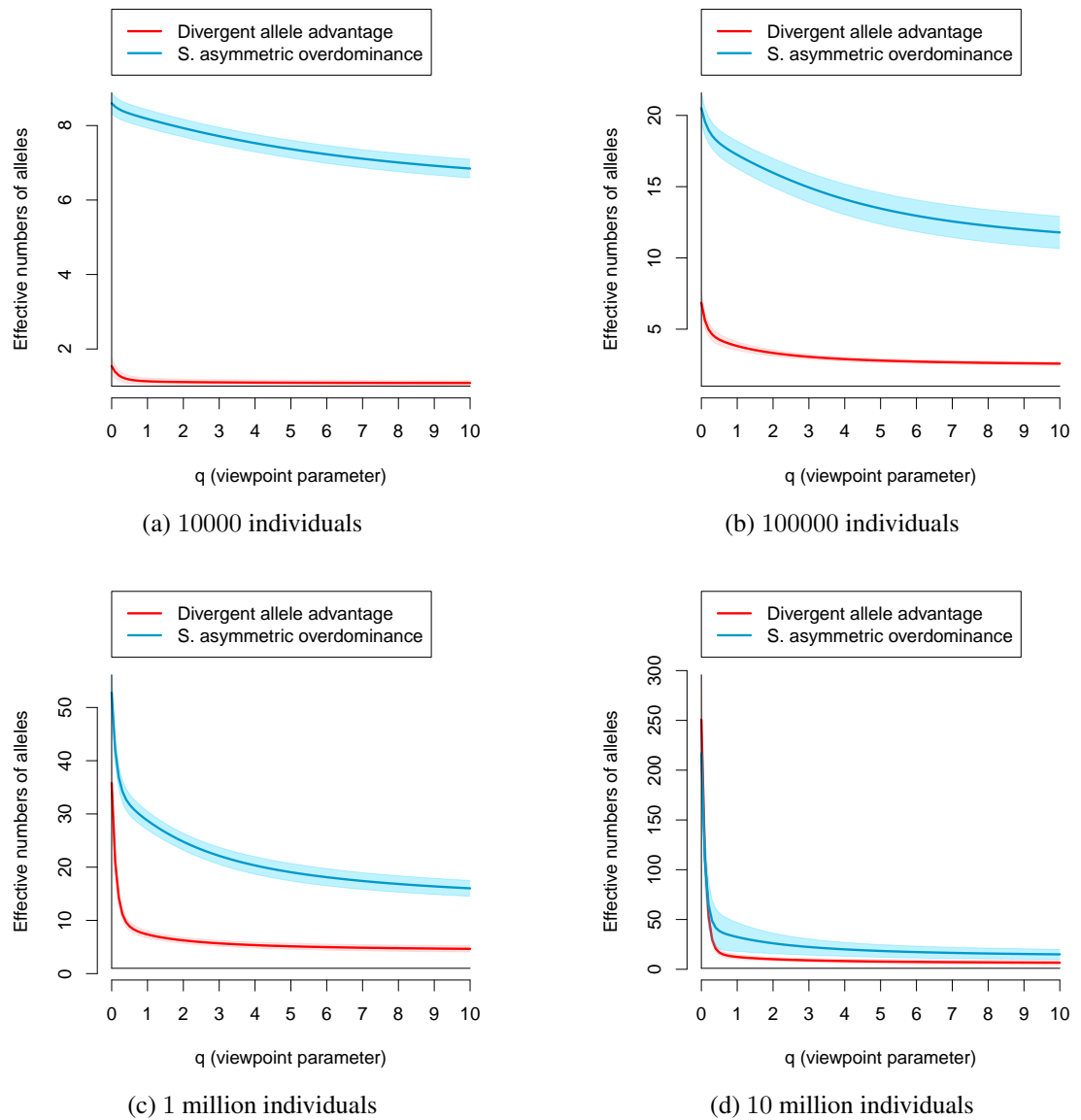


Figure B.35: Diversity profiles (naive diversity) for the DAA and SAsOD models, setting 4 and population sizes of 10000, 100000, 1 million and 10 million.

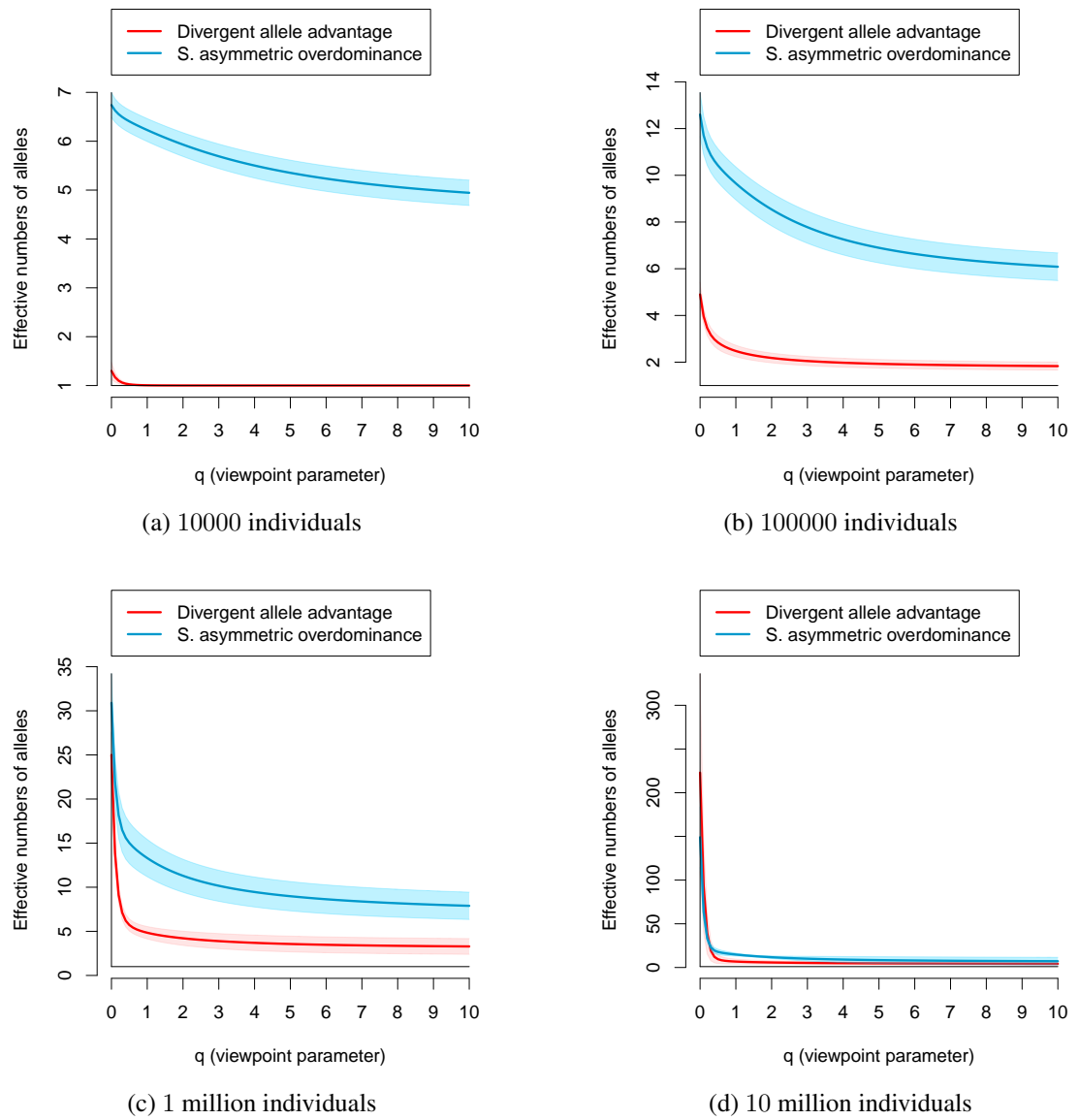


Figure B.36: Diversity profiles (naive diversity) for the DAA and SAsOD models, setting 5 and population sizes of 10000, 100000, 1 million and 10 million.

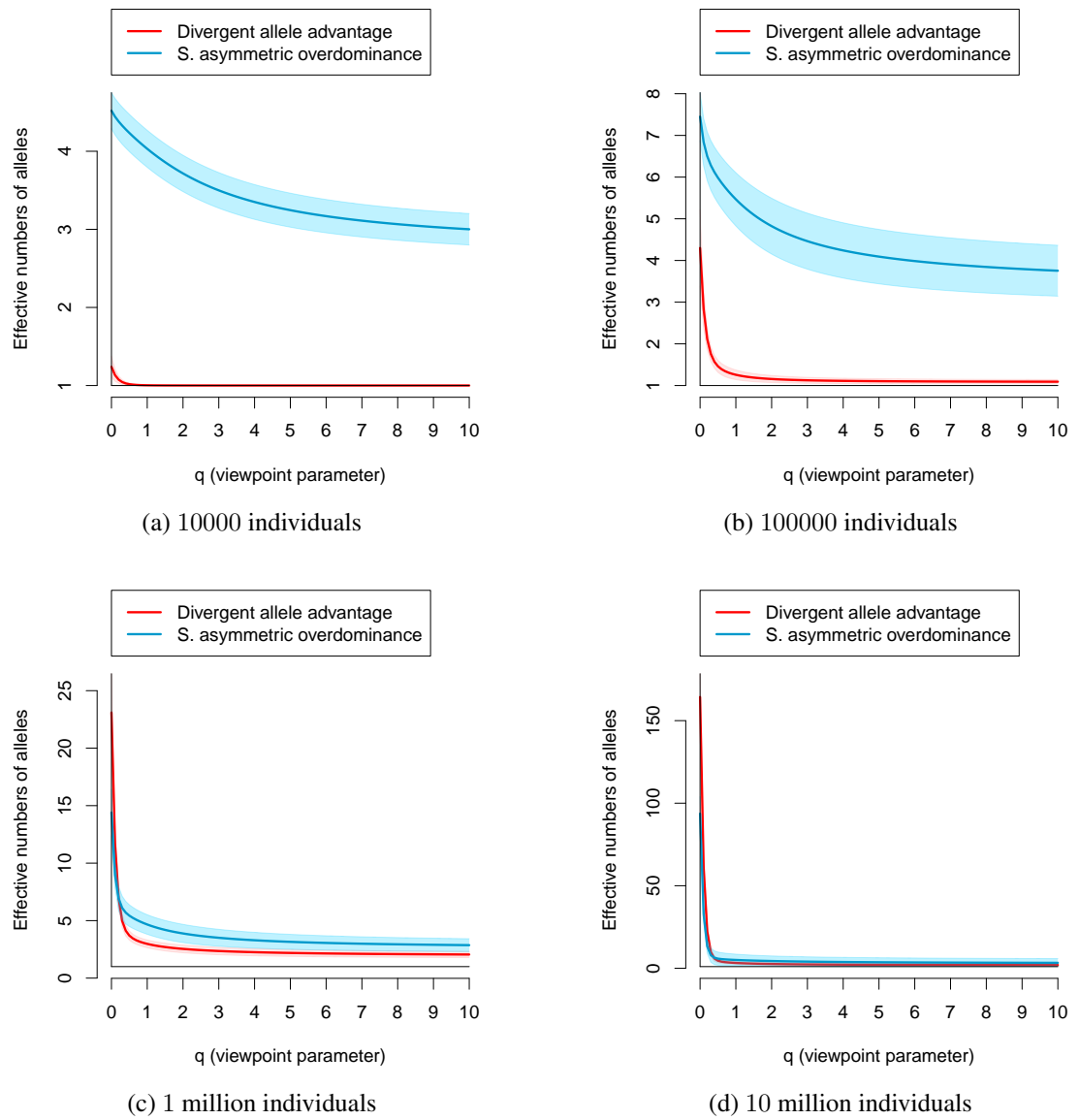


Figure B.37: Diversity profiles (naive diversity) for the DAA and SAsOD models, setting 6 and population sizes of 10000, 100000, 1 million and 10 million.

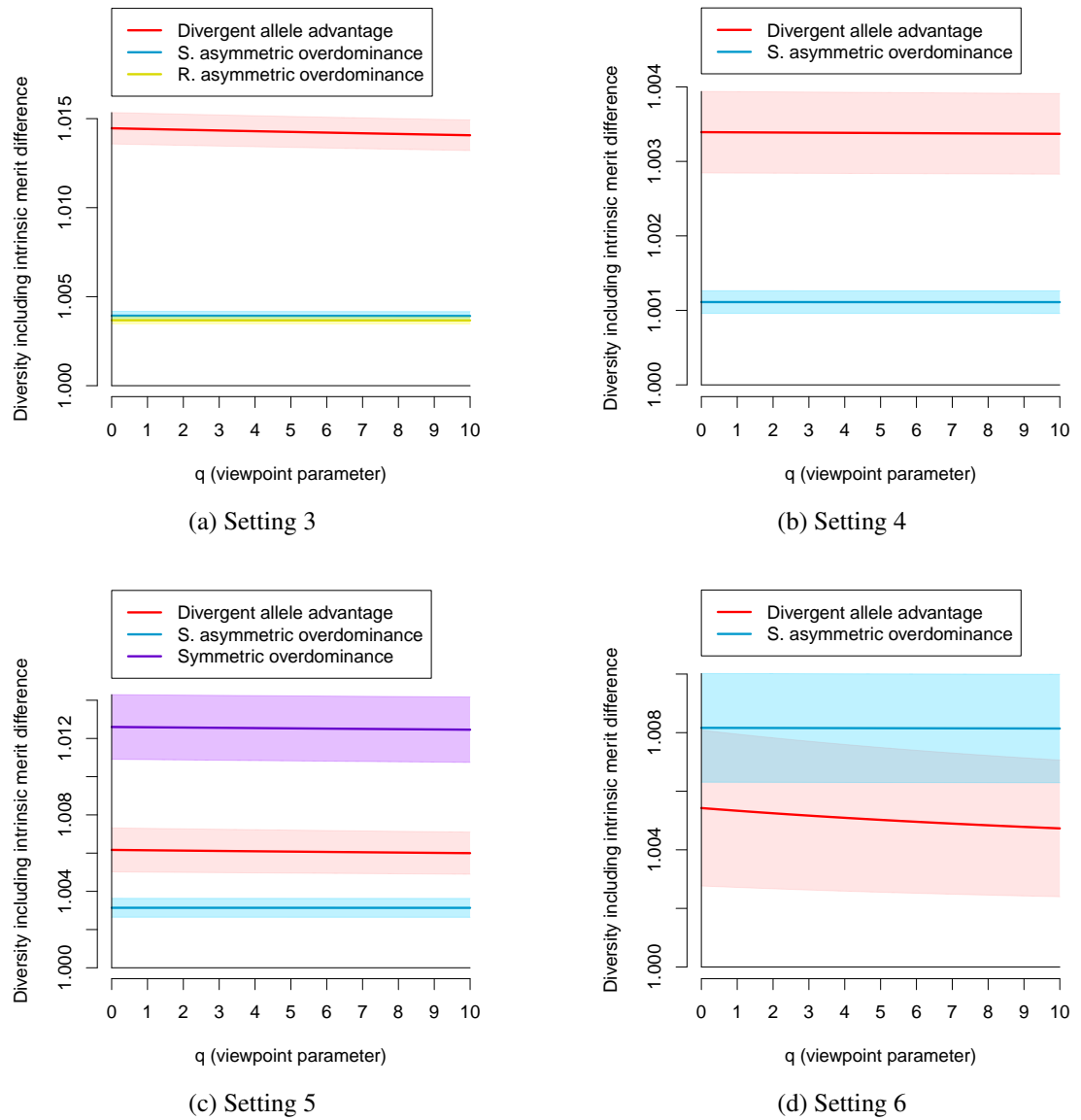


Figure B.38: Diversity profiles when accounting for intrinsic merit differences between alleles for different overdominance models and settings 3, 4, 5 and 6. The population size was 100,000.

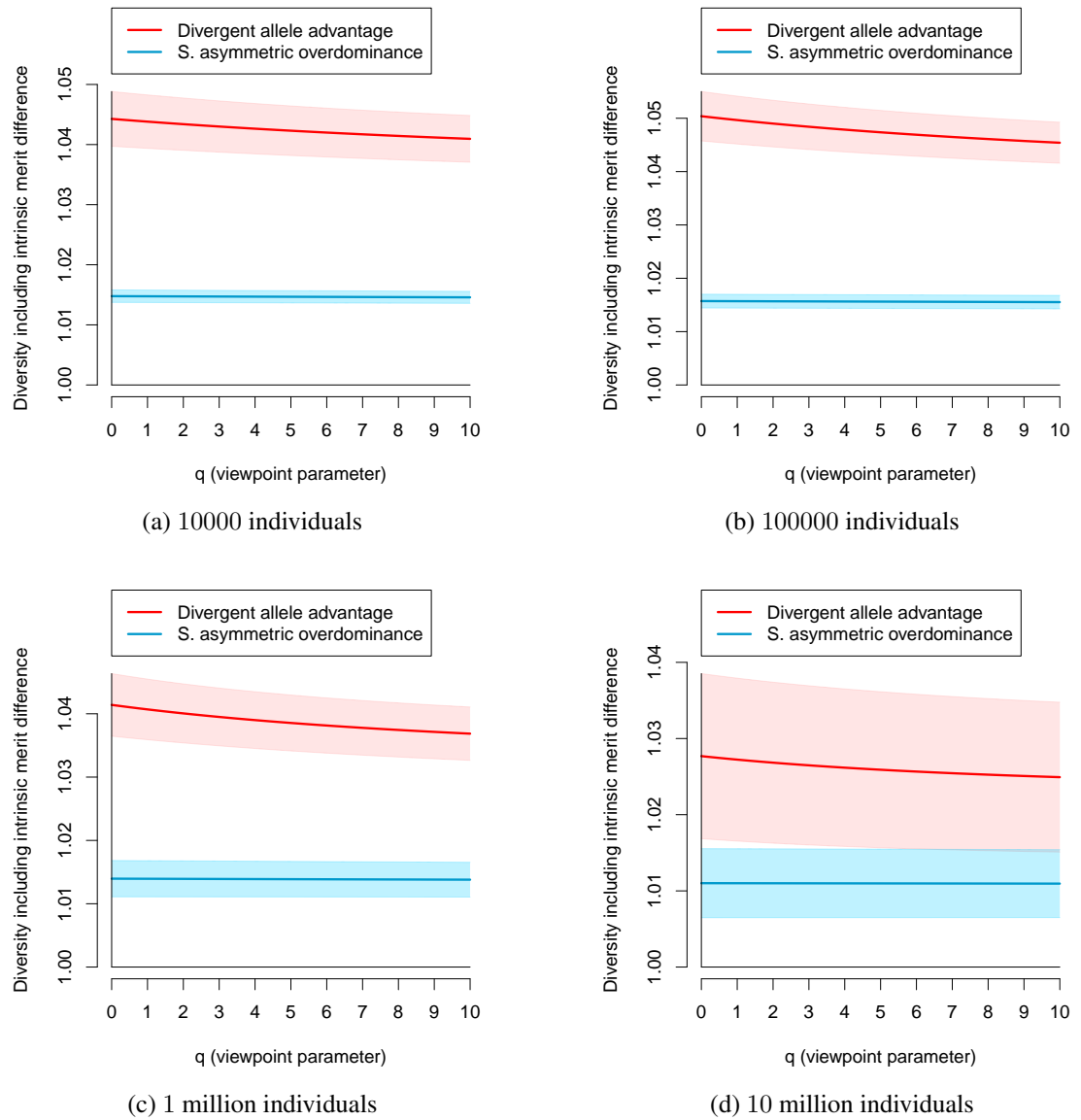


Figure B.39: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and SAsOD models, setting 2 and population sizes of 10000, 100000, 1 million and 10 million.

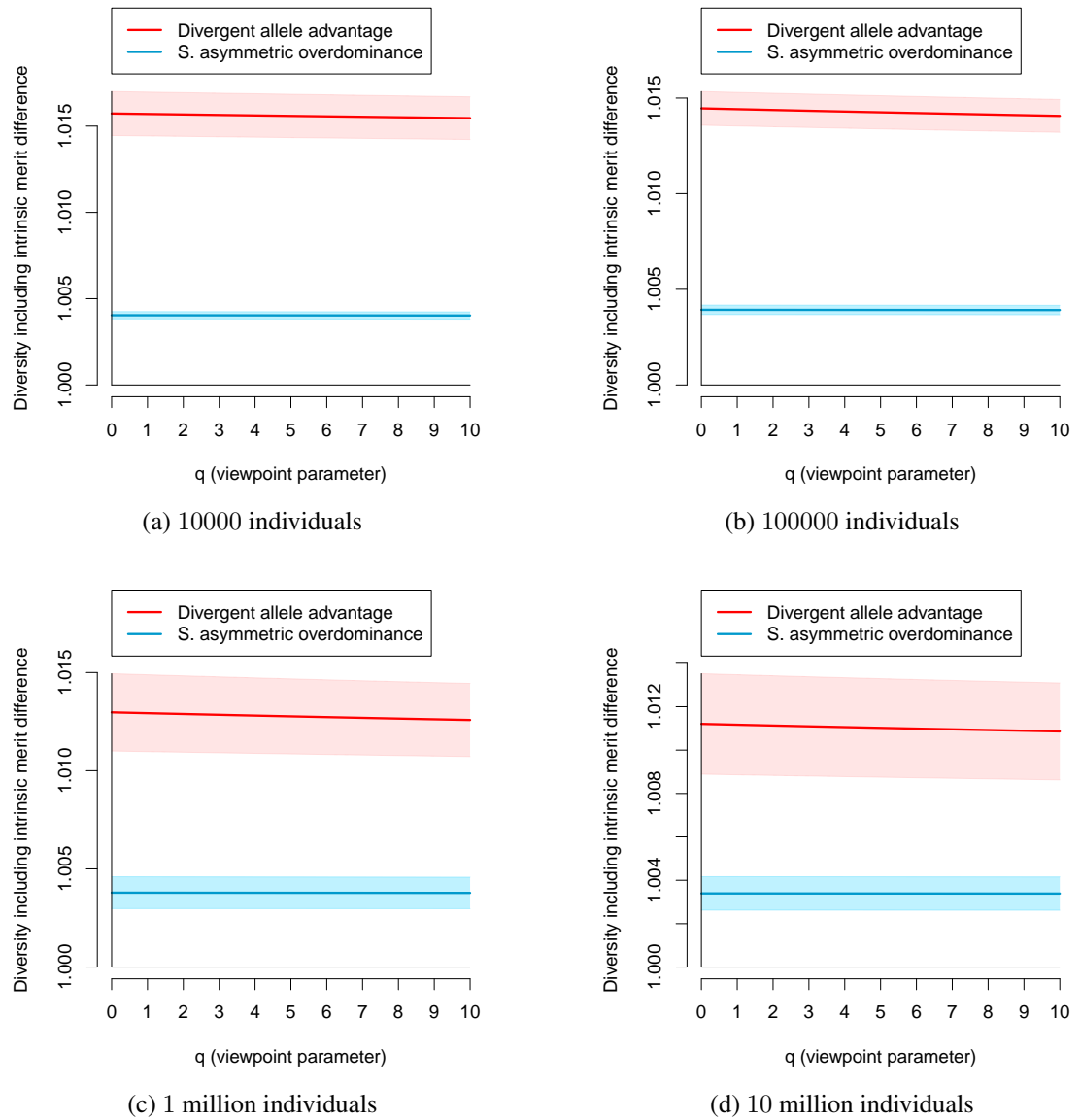


Figure B.40: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and SAsOD models, setting 3 and population sizes of 10000, 100000, 1 million and 10 million.

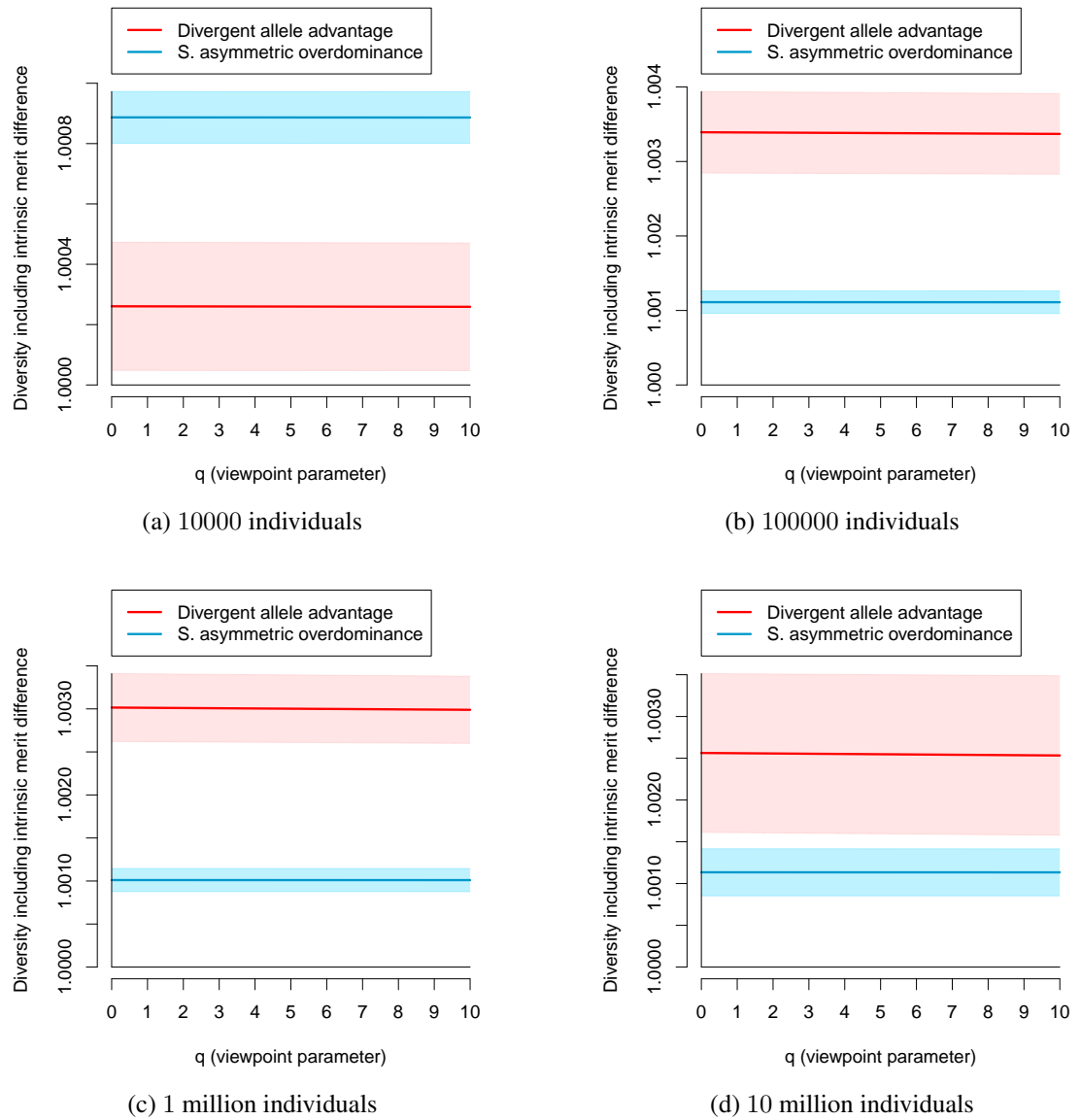


Figure B.41: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and SAsOD models, setting 4 and population sizes of 10000, 100000, 1 million and 10 million.

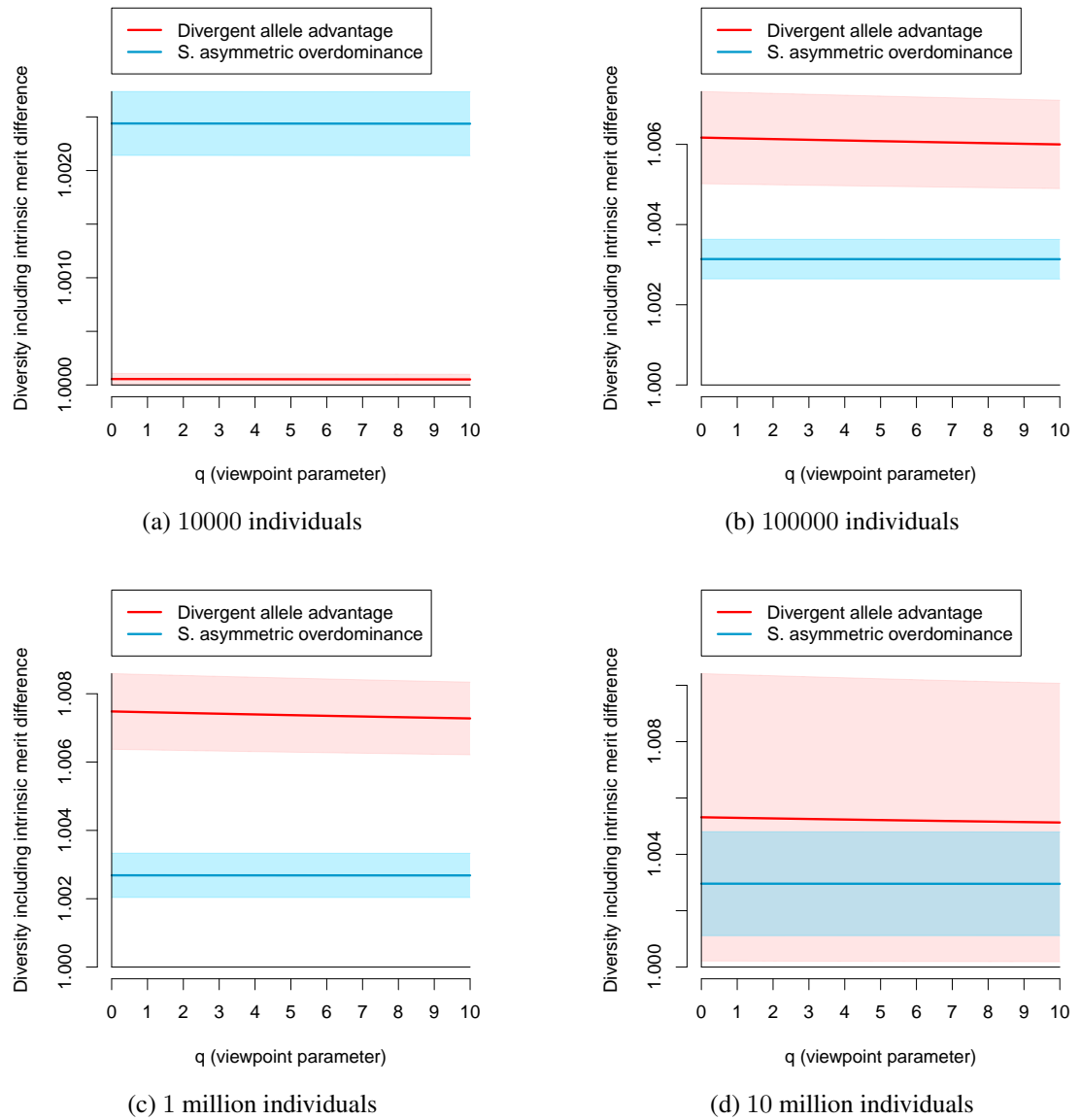


Figure B.42: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and SAsOD models, setting 5 and population sizes of 10000, 100000, 1 million and 10 million.

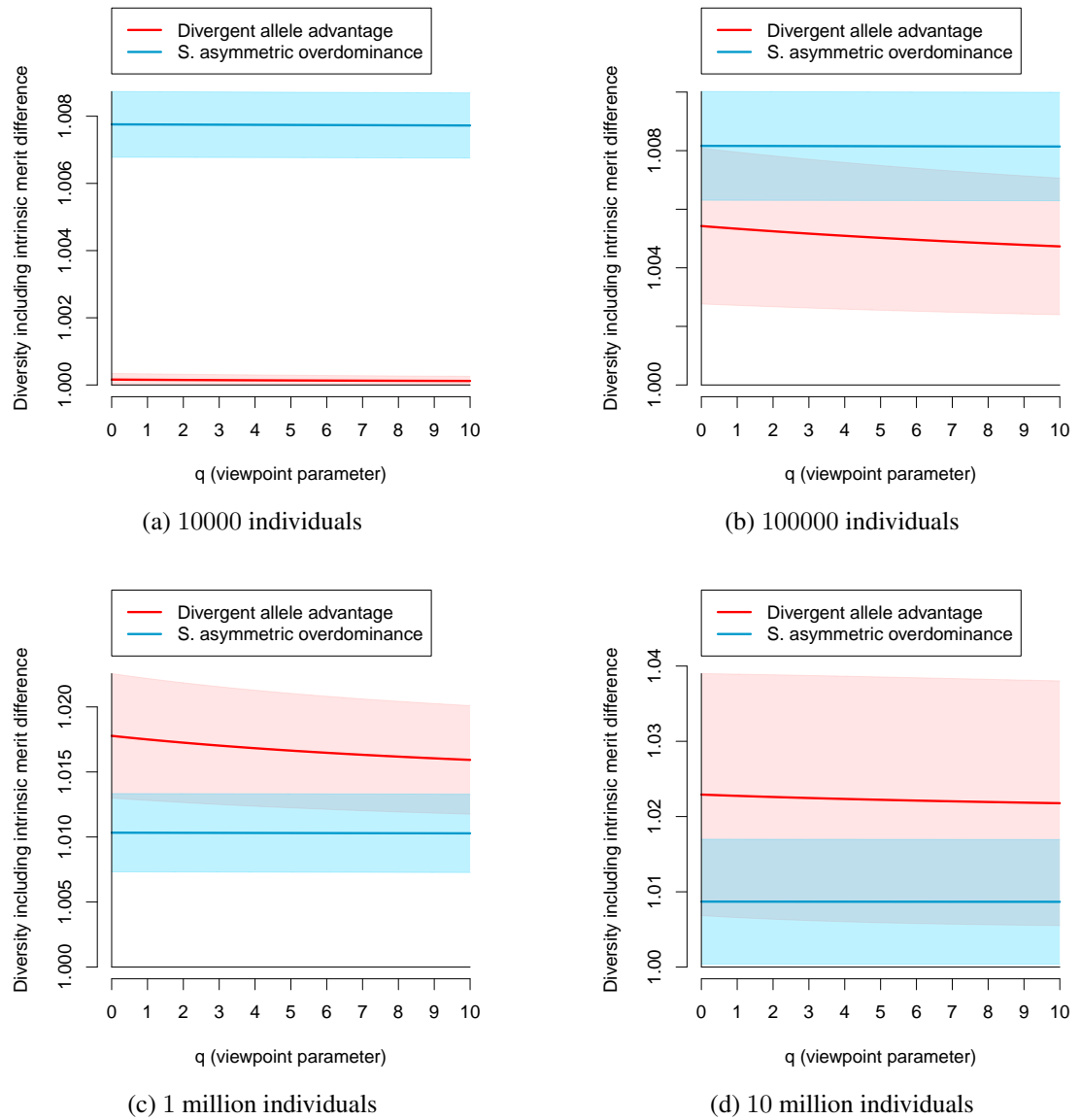
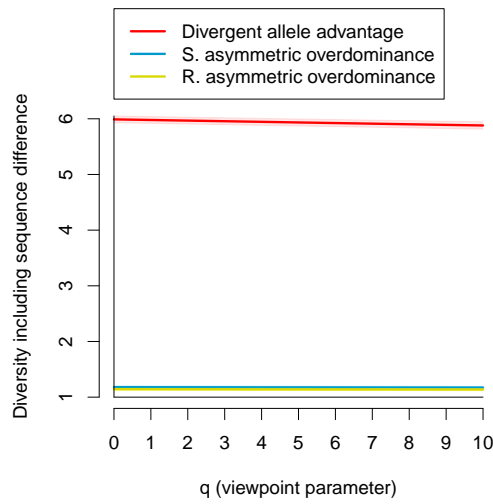
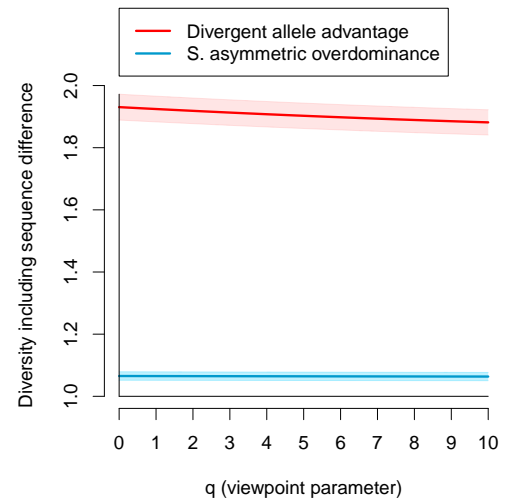


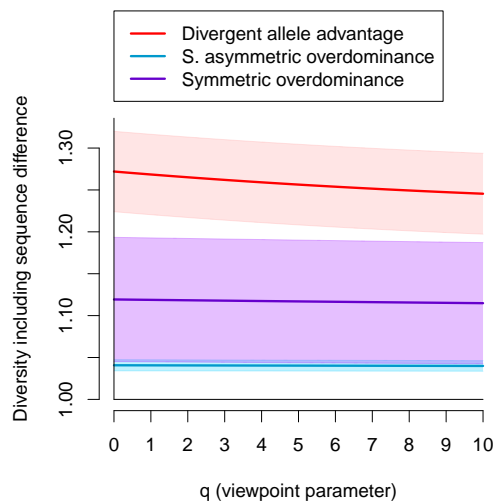
Figure B.43: Diversity profiles when accounting for intrinsic merit differences between alleles for the DAA and SAsOD models, setting θ and population sizes of 10000, 100000, 1 million and 10 million.



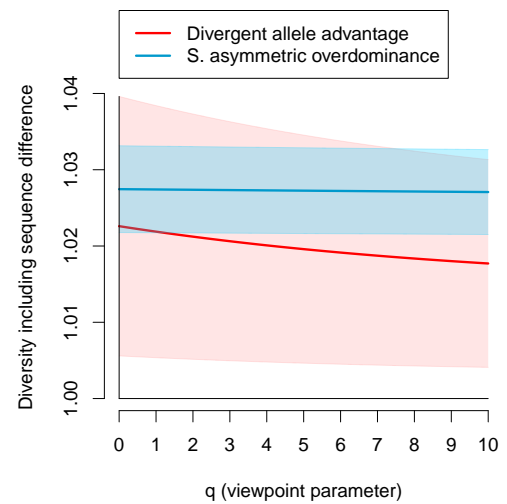
(a) Setting 3



(b) Setting 4



(c) Setting 5



(d) Setting 6

Figure B.44: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for different overdominance models and settings 3, 4, 5 and 6. The population size was 100000.

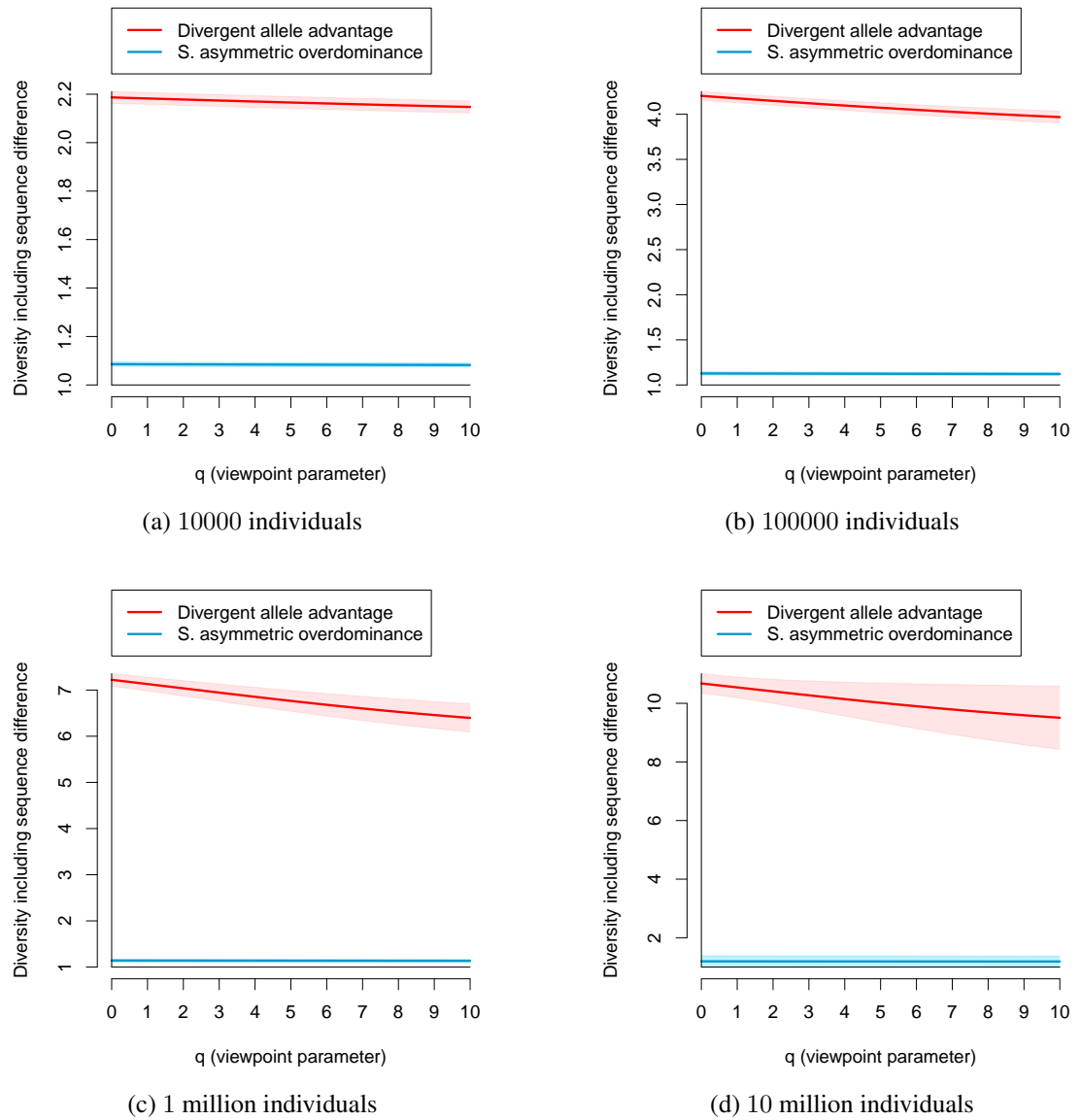


Figure B.45: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and SAsOD models, setting 2 and population sizes of 10000, 100000, 1 million and 10 million.

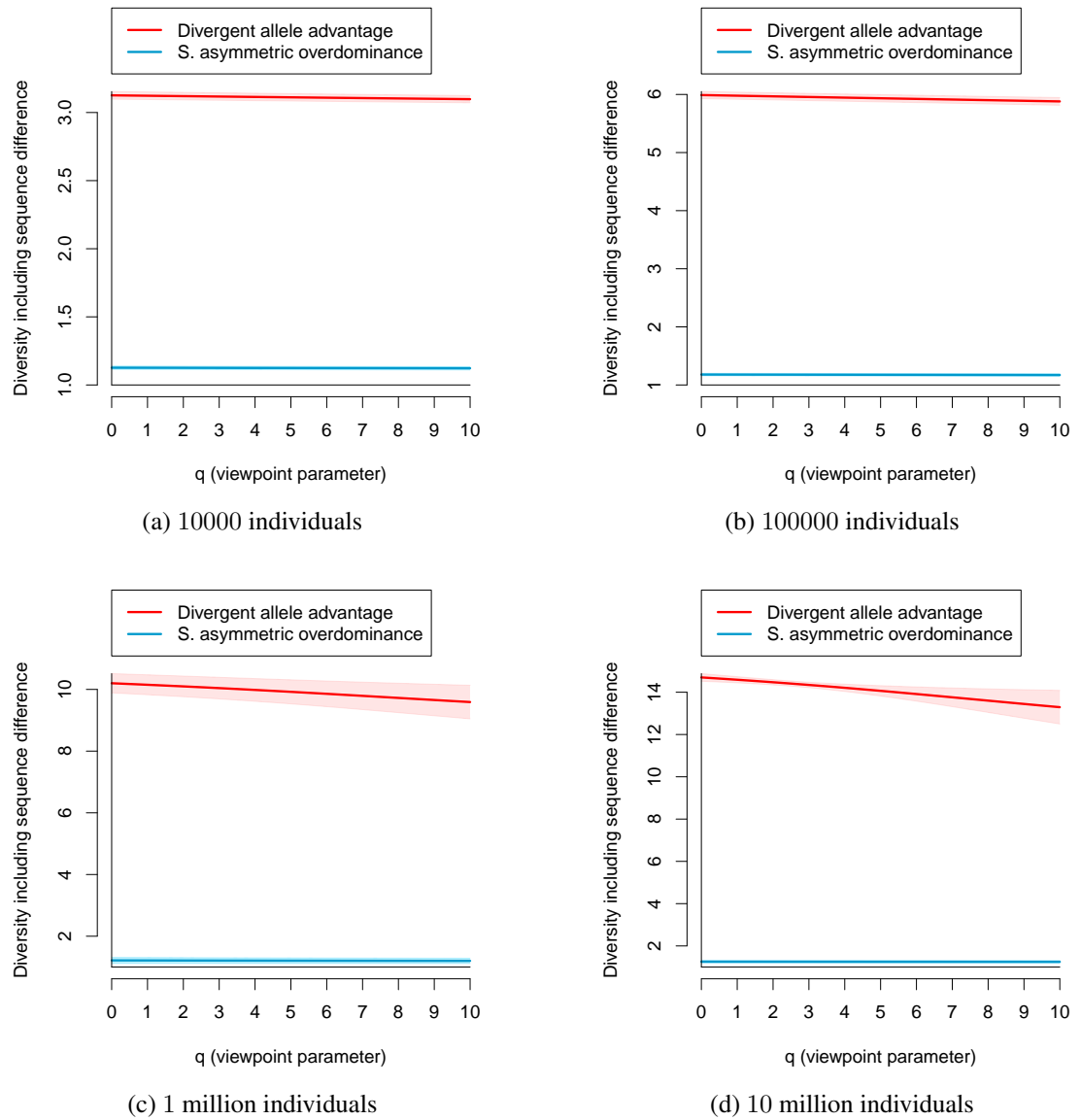


Figure B.46: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and SAsOD models, setting 3 and population sizes of 10000, 100000, 1 million and 10 million.

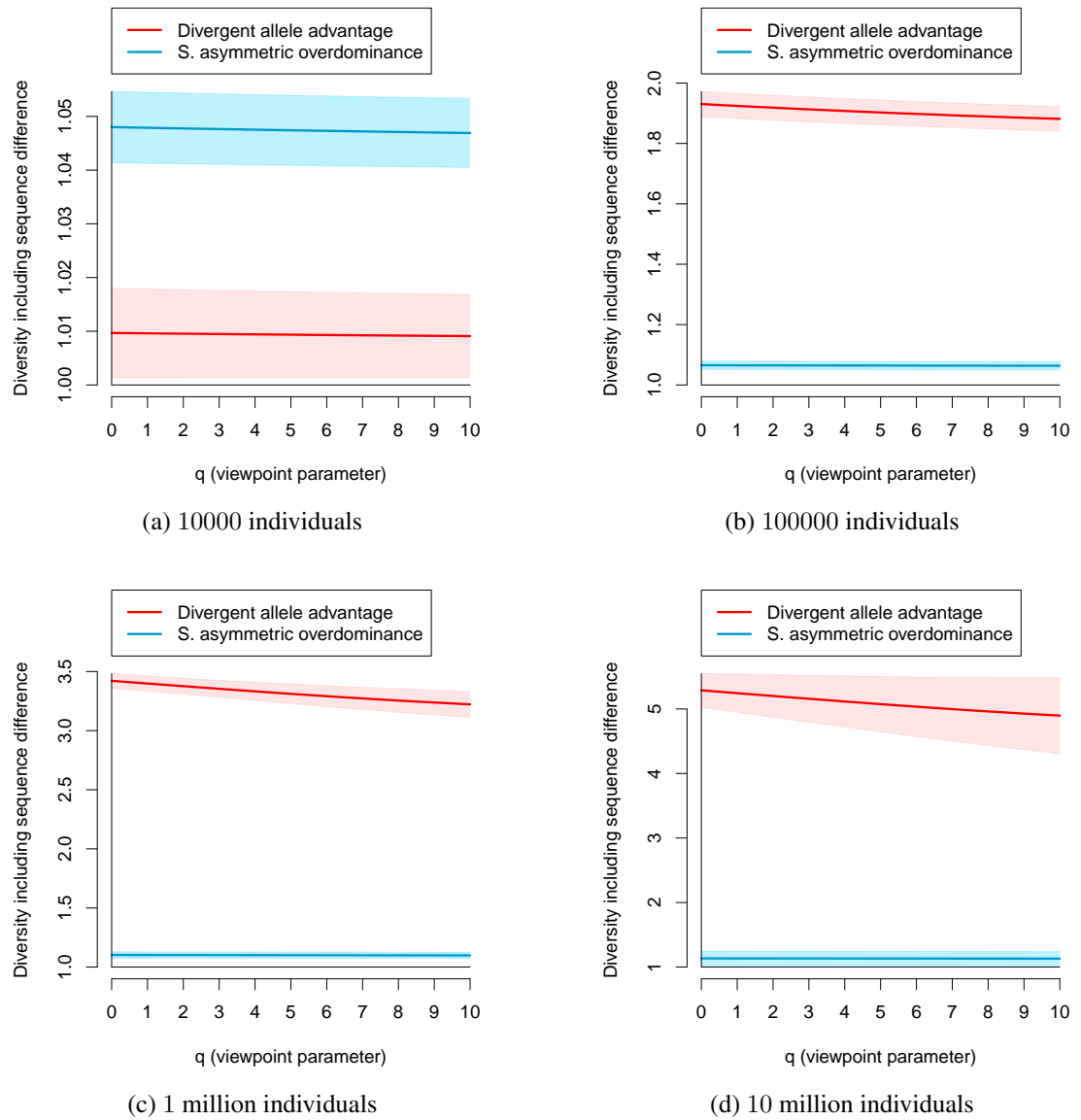


Figure B.47: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and SAsOD models, setting 4 and population sizes of 10000, 100000, 1 million and 10 million.

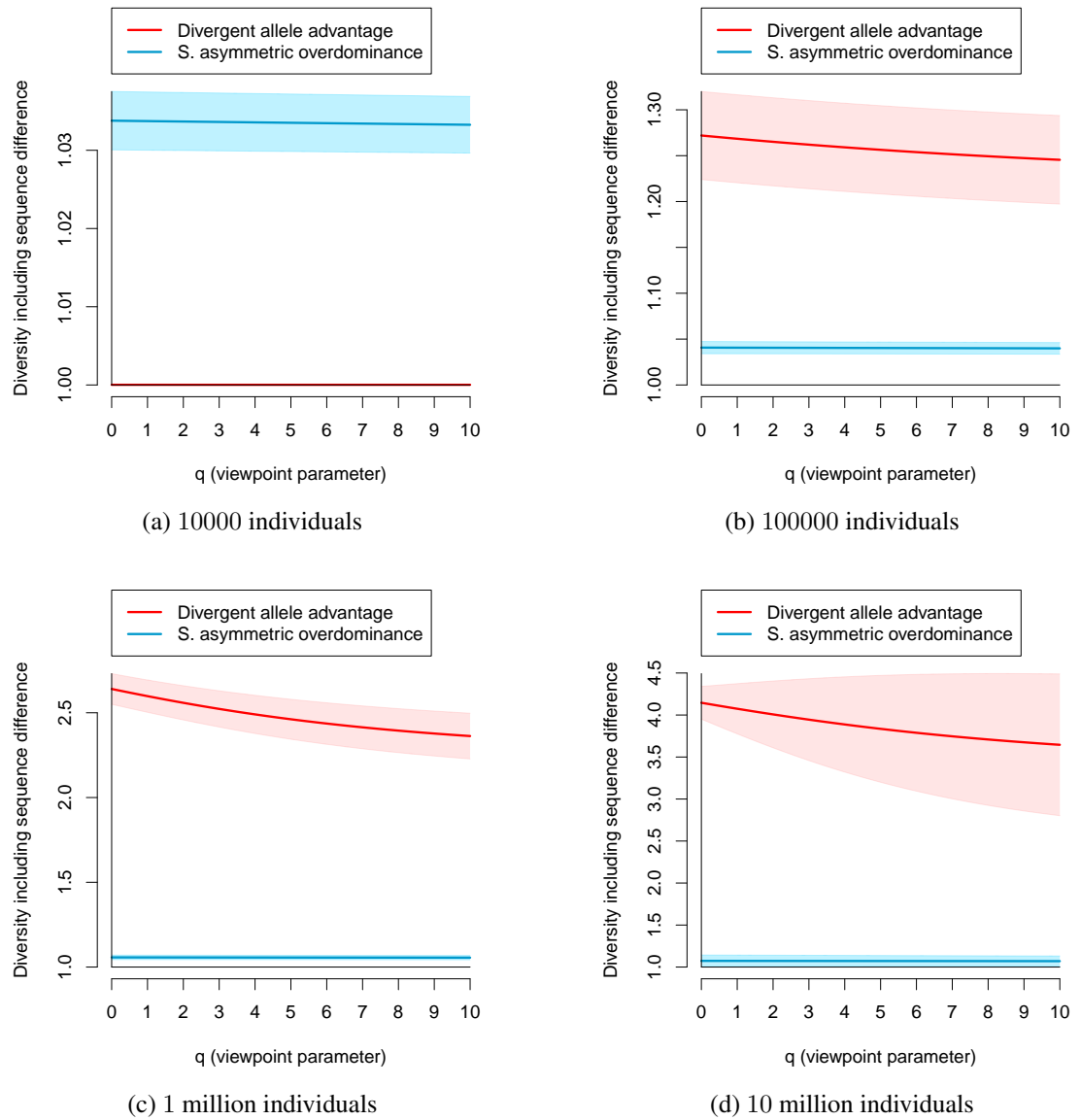


Figure B.48: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and SAsOD models, setting 5 and population sizes of 10000, 100000, 1 million and 10 million.

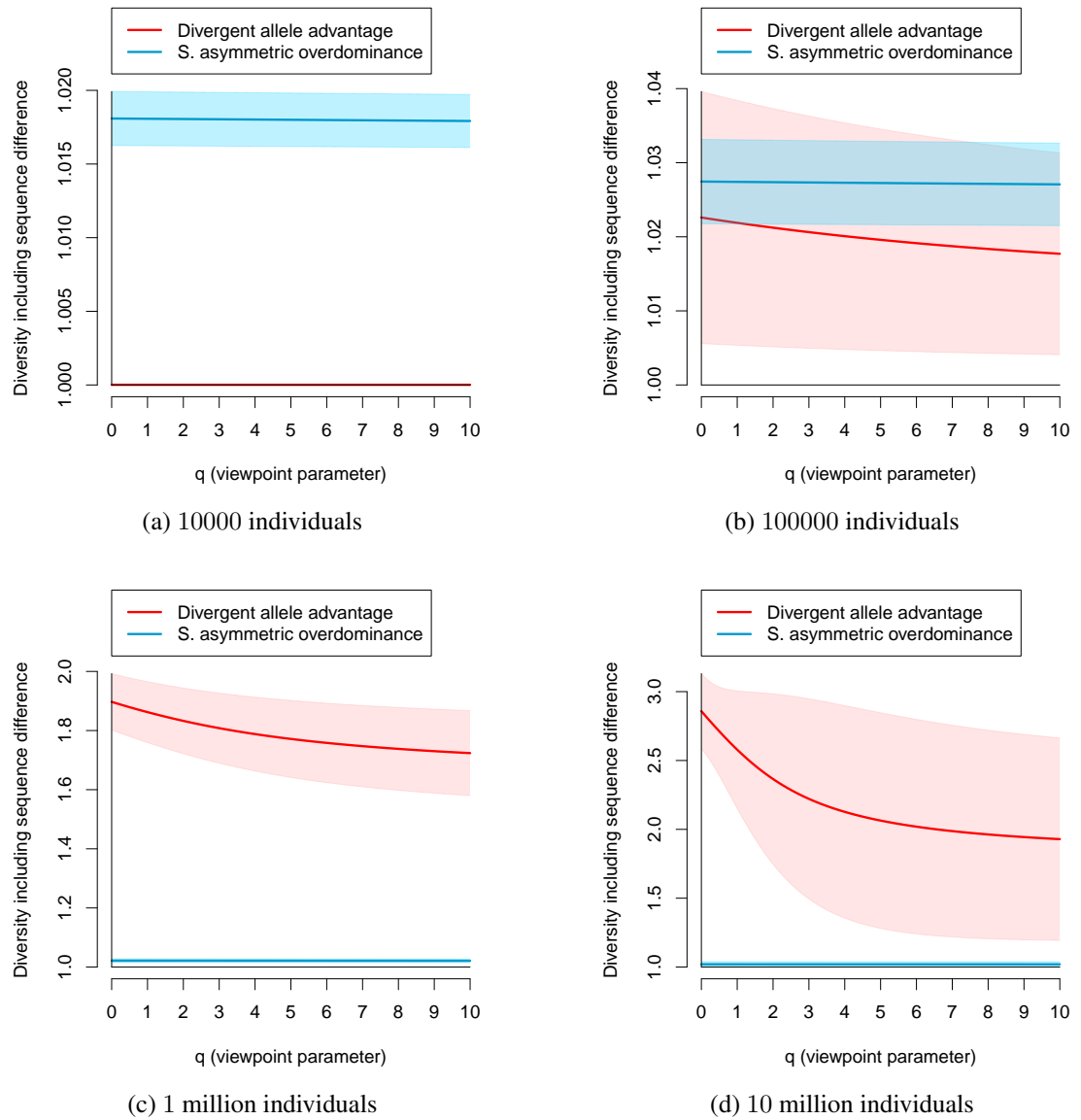


Figure B.49: Diversity profiles when accounting for differences in the amino acid sequences encoded by the alleles for the DAA and SAsOD models, setting θ and population sizes of 10000, 100000, 1 million and 10 million.

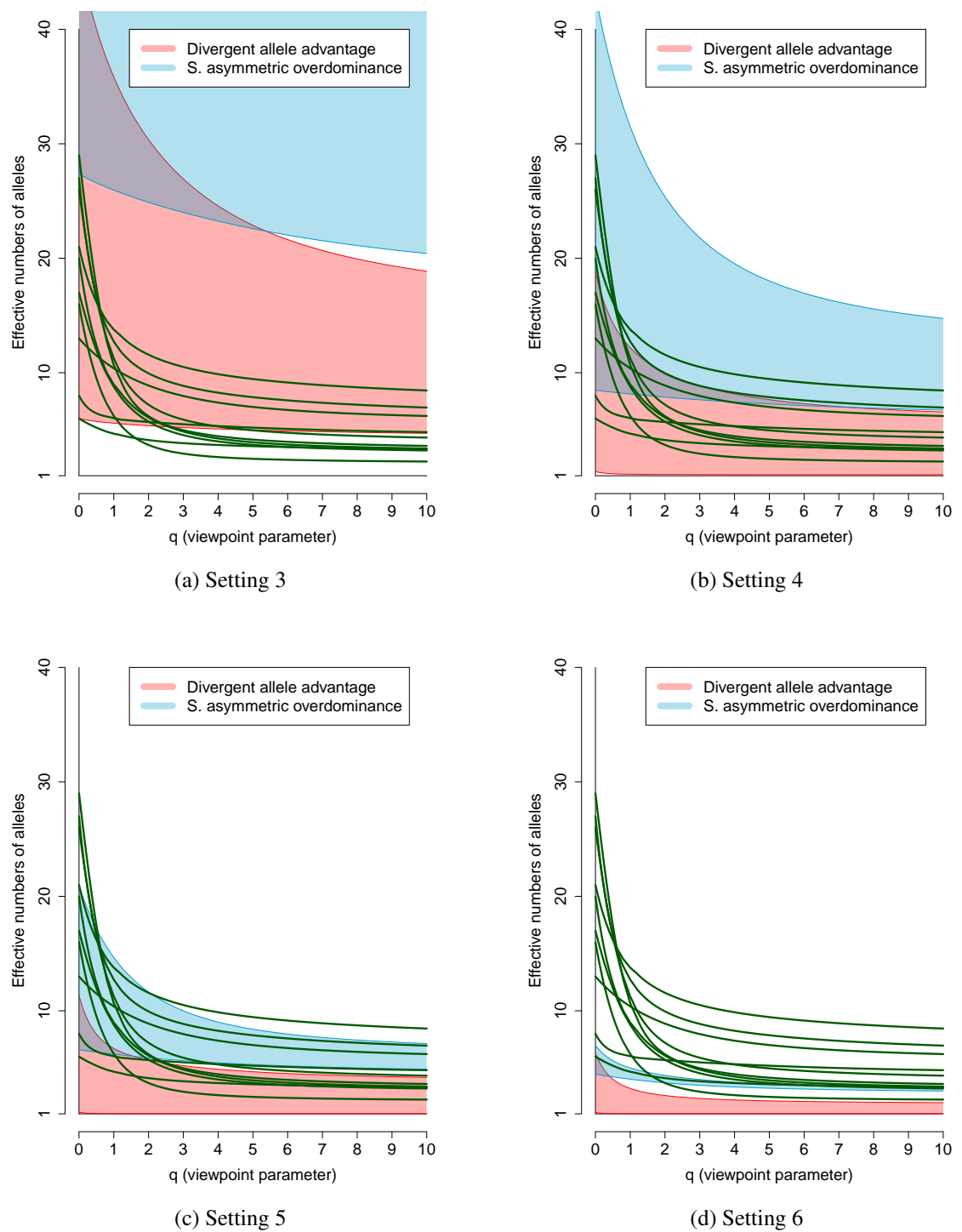
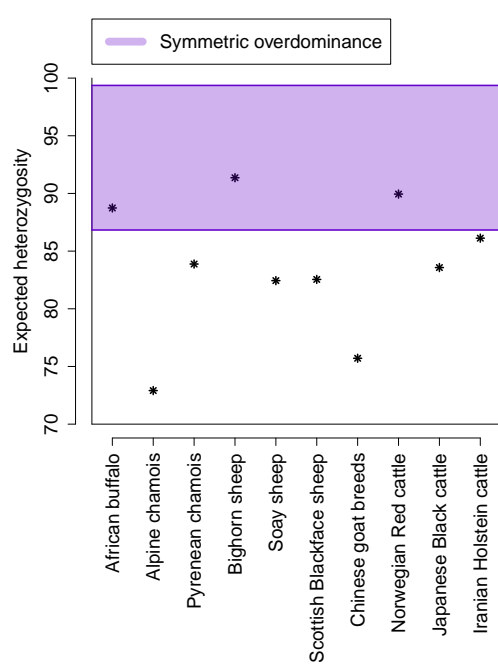
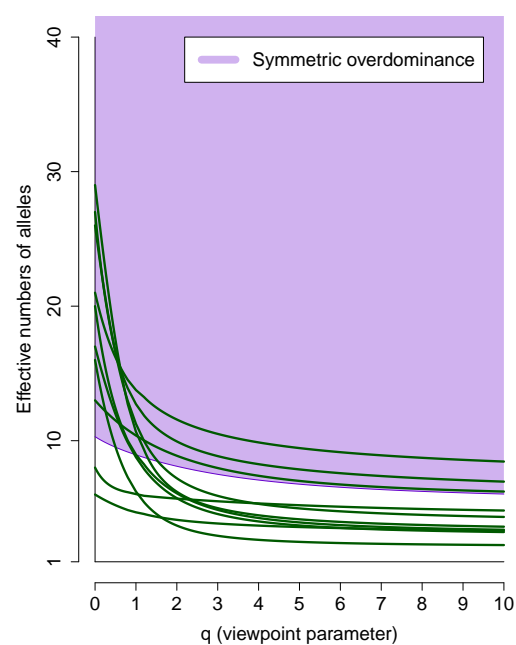


Figure B.50: Diversity calculated from observed MHC allele frequencies and samples drawn from the final gene pools of alternative overdominance models for settings 3, 4, 5 and 6.



(a) Setting 5



(b) Setting 5

Figure B.51: Expected heterozygosity and diversity profile for observed MHC allele frequencies and samples drawn from the final gene pool of the SOD model.

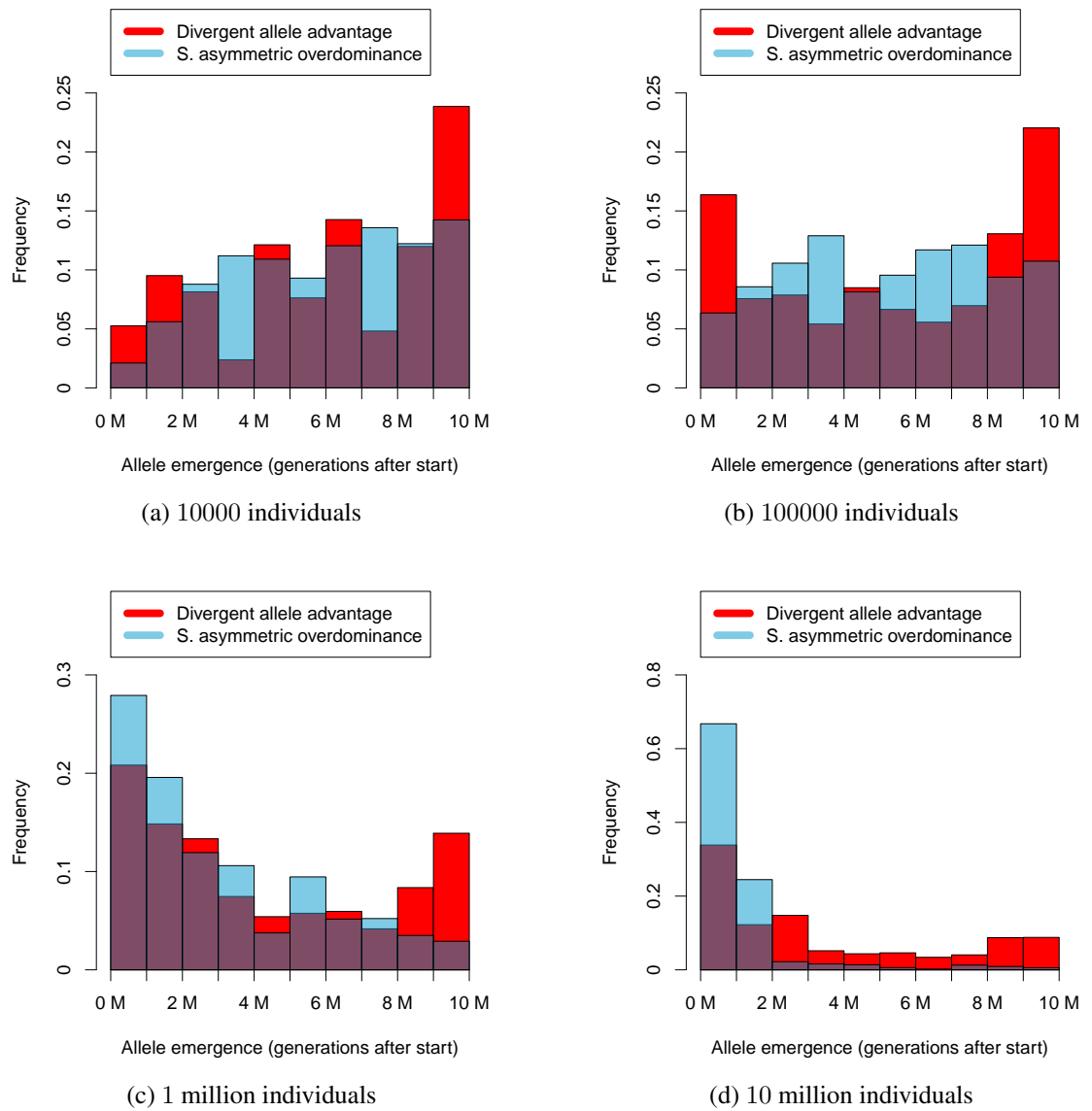


Figure B.52: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 1. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

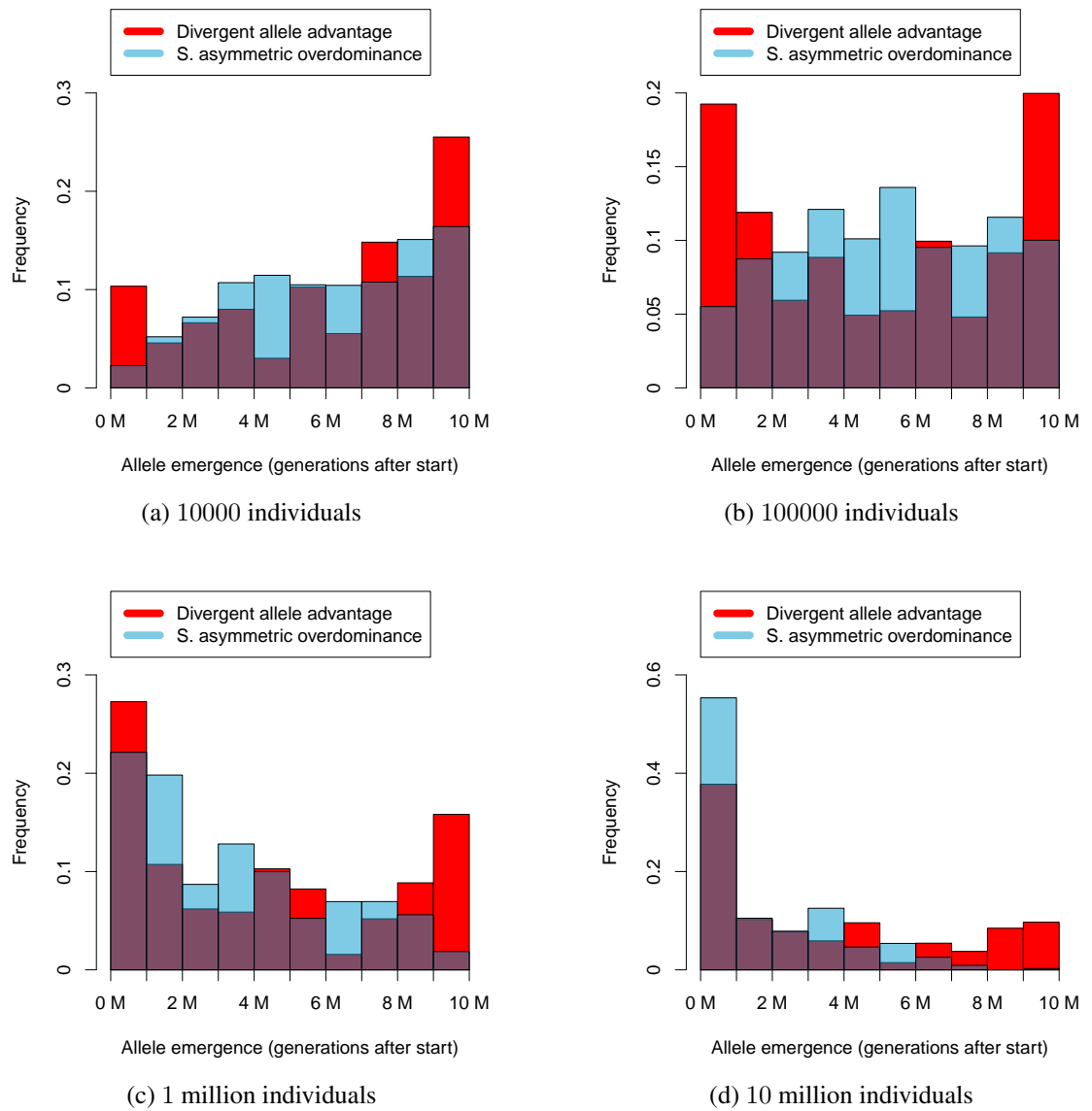


Figure B.53: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 2. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

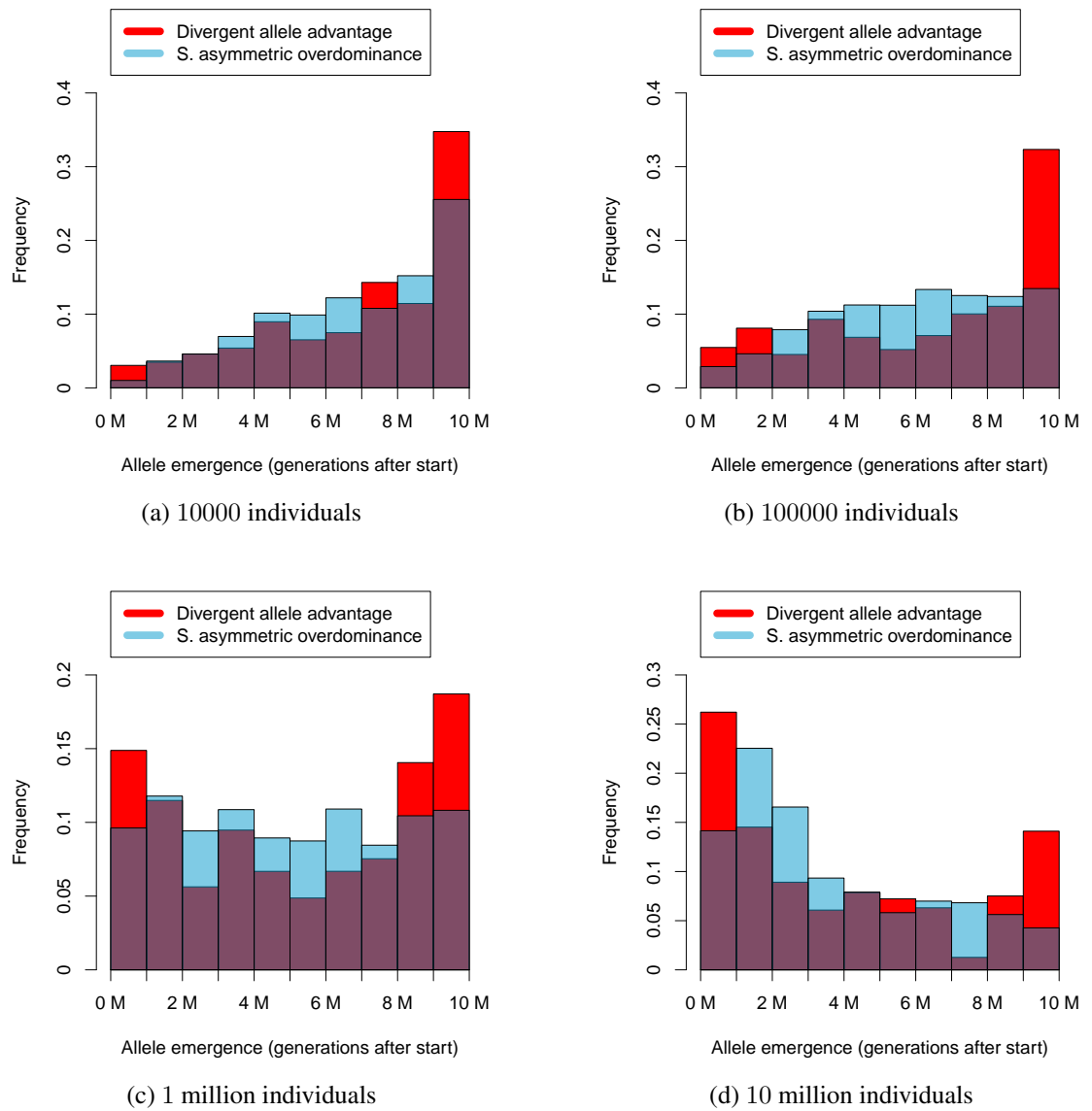


Figure B.54: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 3. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

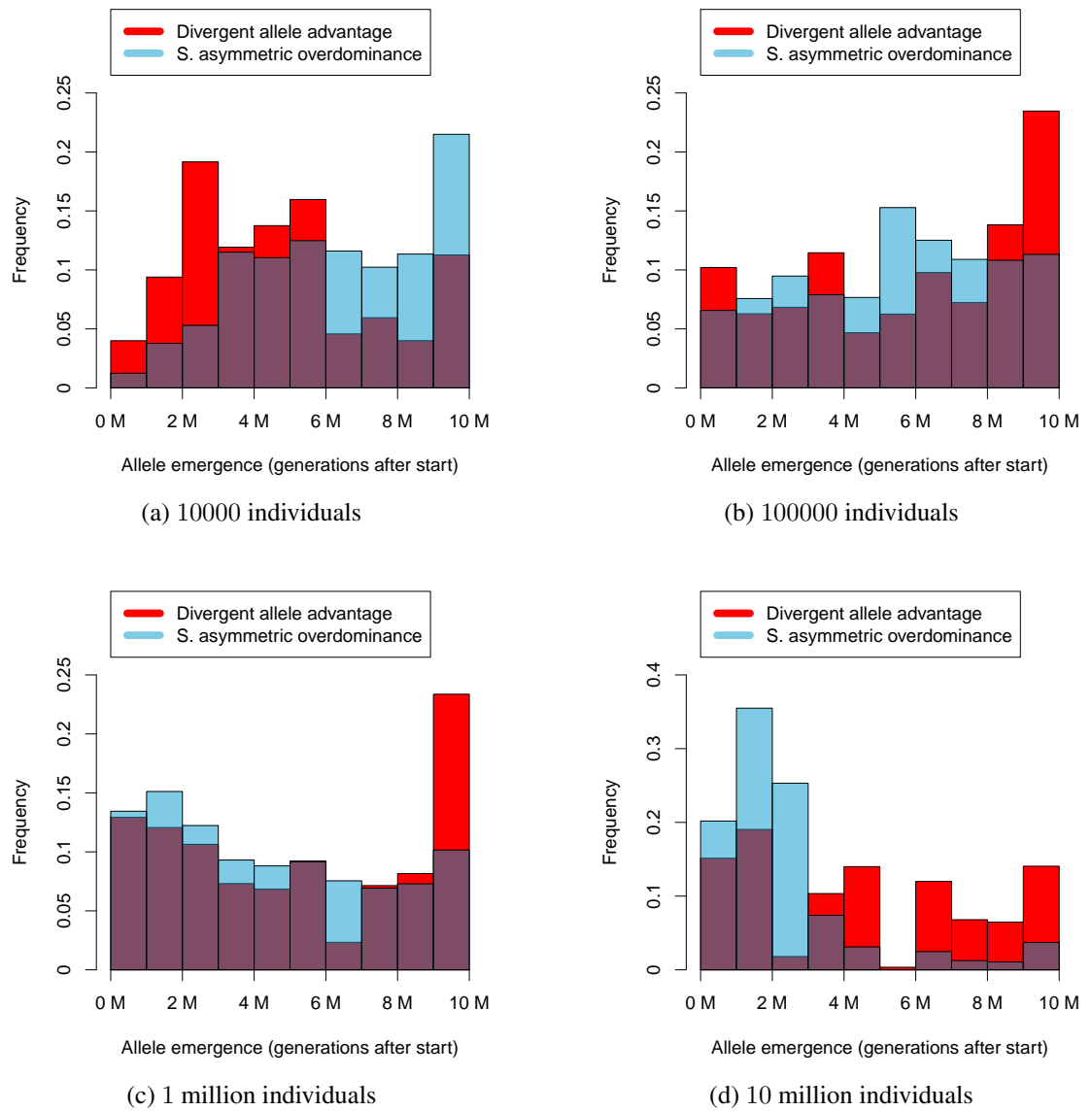


Figure B.55: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 4. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

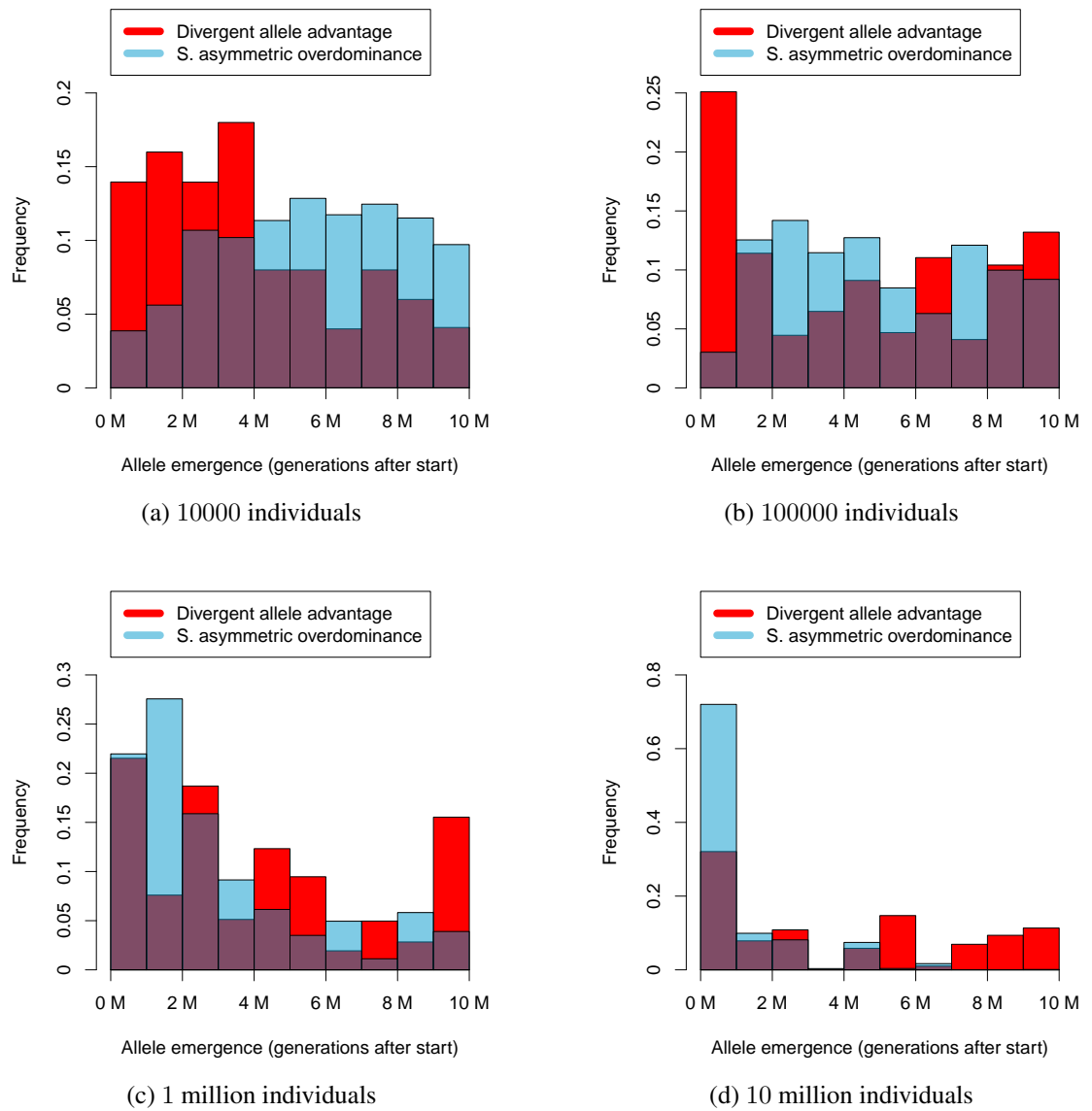


Figure B.56: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 5. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

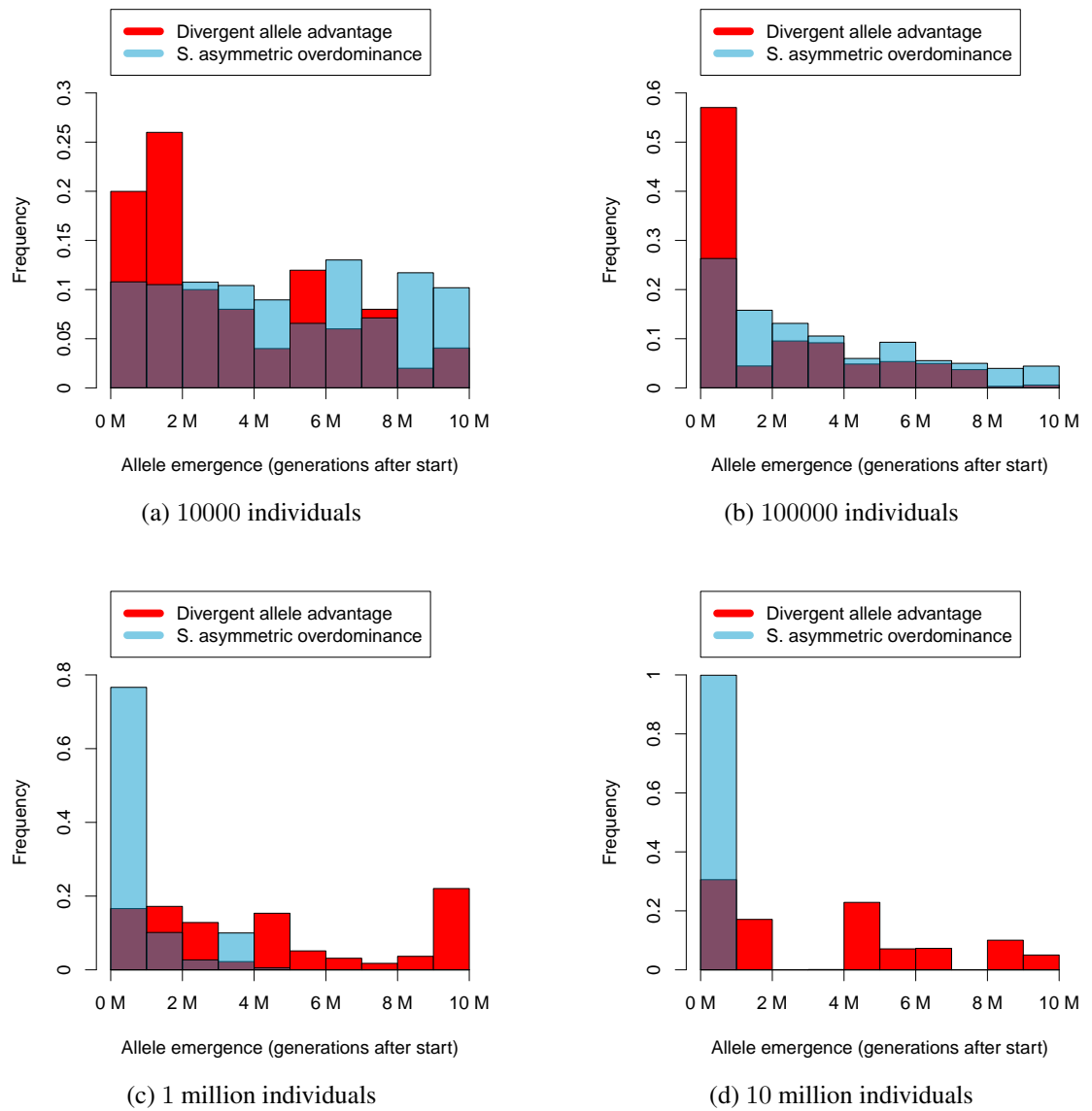
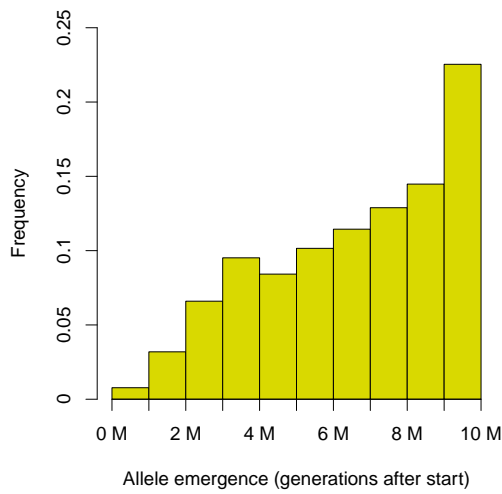
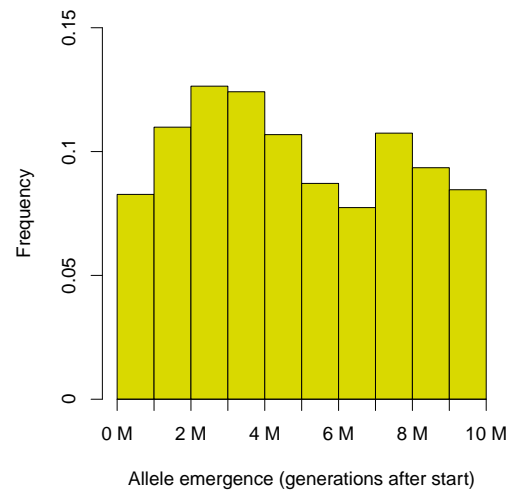


Figure B.57: Distribution of emergence times of alleles for the DAA and SAsOD models and setting 6. The population size is 10000, 100000, 1 million and 10 million (from top left to bottom right).

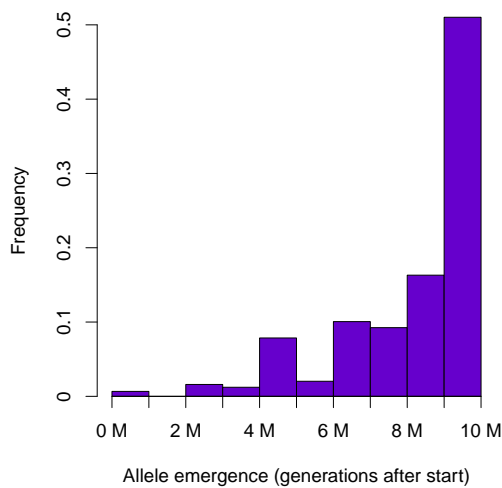


(a) RAsOD model, setting 3, 10000 individuals

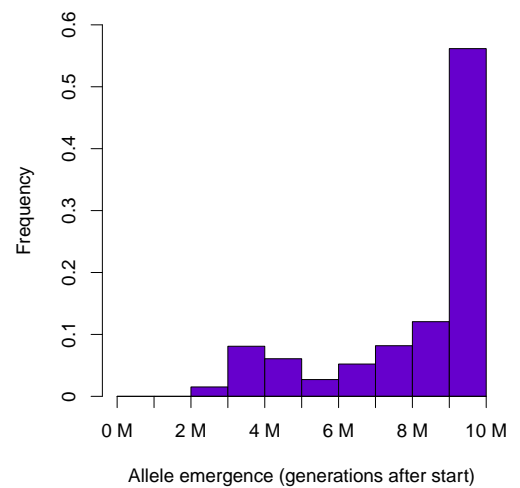


(b) RAsOD model, setting 3, 1 million individuals

Figure B.58: Distribution of emergence times of alleles for the RAsOD model and setting 3. The population size is 10000 (left) and 1 million (right), respectively.



(a) SOD model, setting 5, 10000 individuals



(b) SOD model, setting 5, 100000 individuals

Figure B.59: Distribution of emergence times of alleles for the SOD model and setting 5. The population size is 10000 (left) and 100000 (right), respectively.

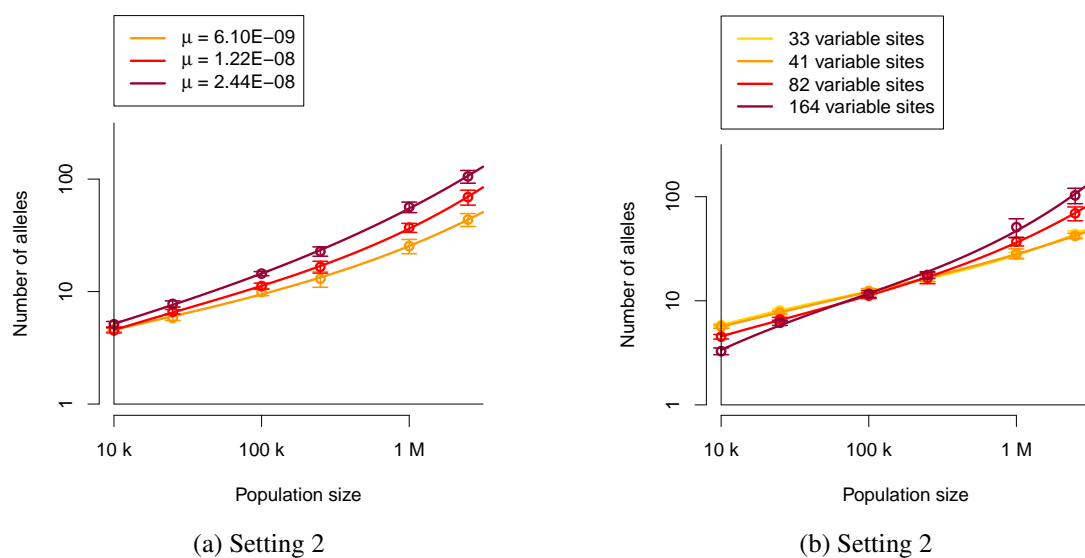
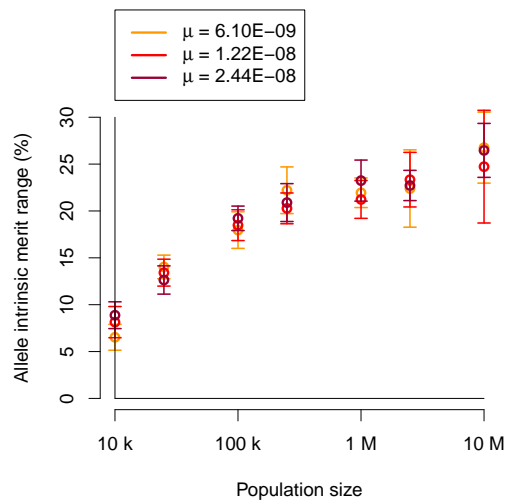
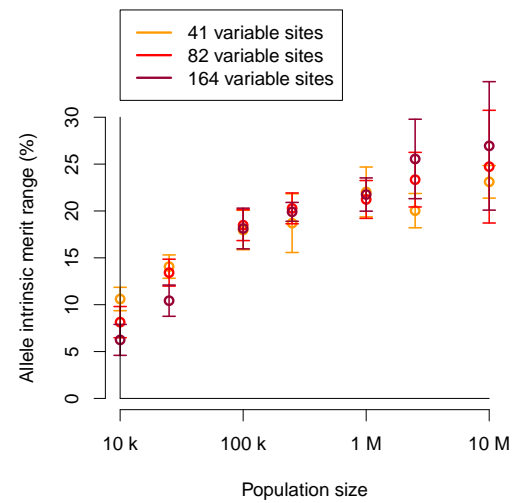


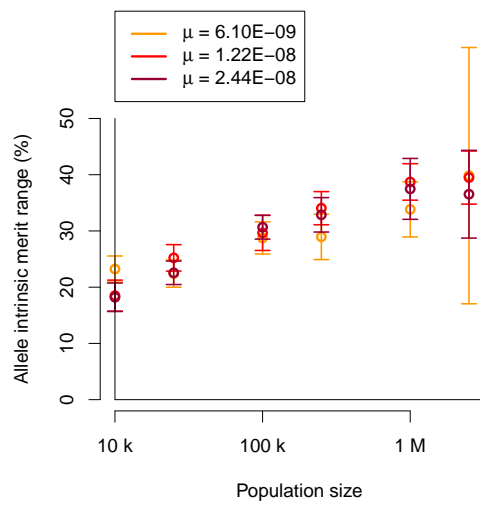
Figure B.60: Number of alleles vs. population size for setting 2 and different mutation rates (left) and number of variable sites (right).



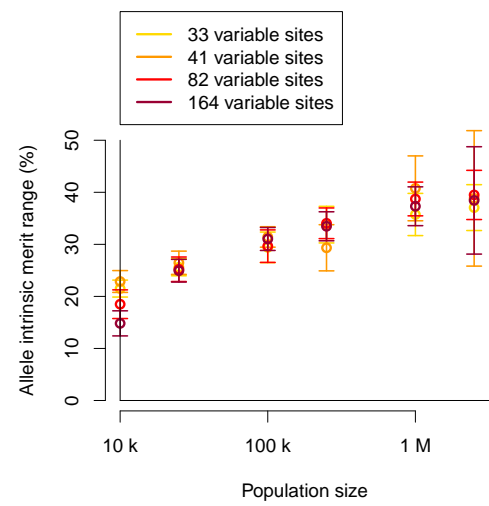
(a) Setting 1



(b) Setting 1

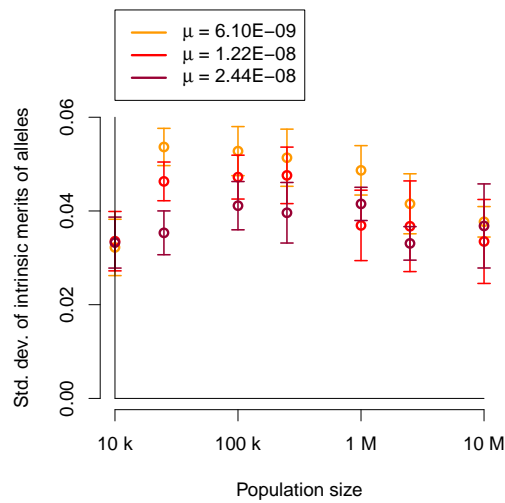


(c) Setting 2

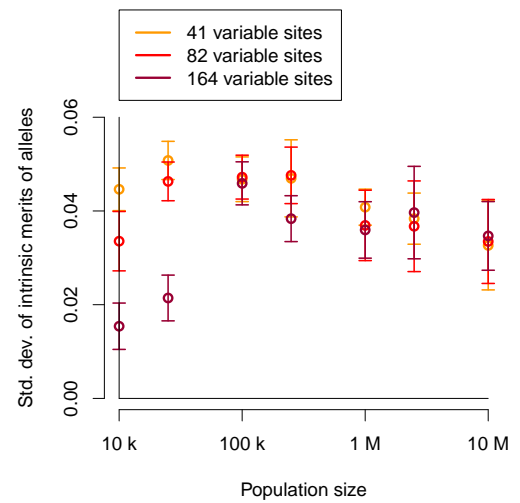


(d) Setting 2

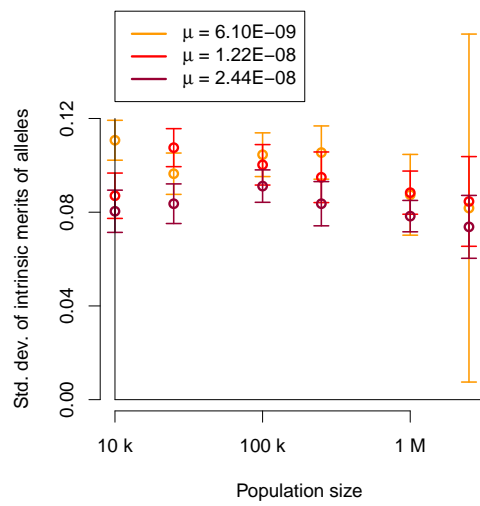
Figure B.61: Range of allele intrinsic merits vs. population size for settings 1 and 2 and different mutation rates (left) and number of variable sites (right).



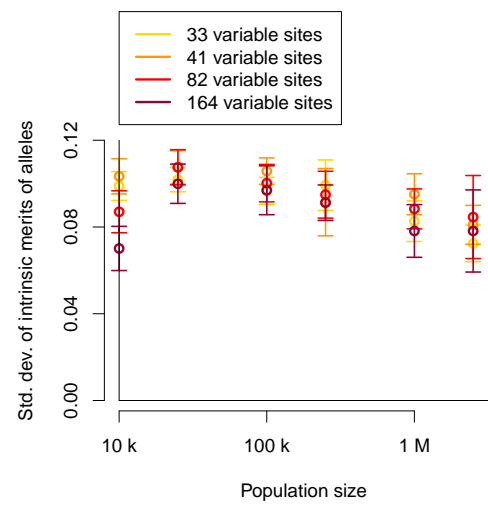
(a) Setting 1



(b) Setting 1



(c) Setting 2



(d) Setting 2

Figure B.62: Weighted standard deviation in allele intrinsic merit vs. population size for settings 1 and 2 and different mutation rates (left) and number of variable sites (right).

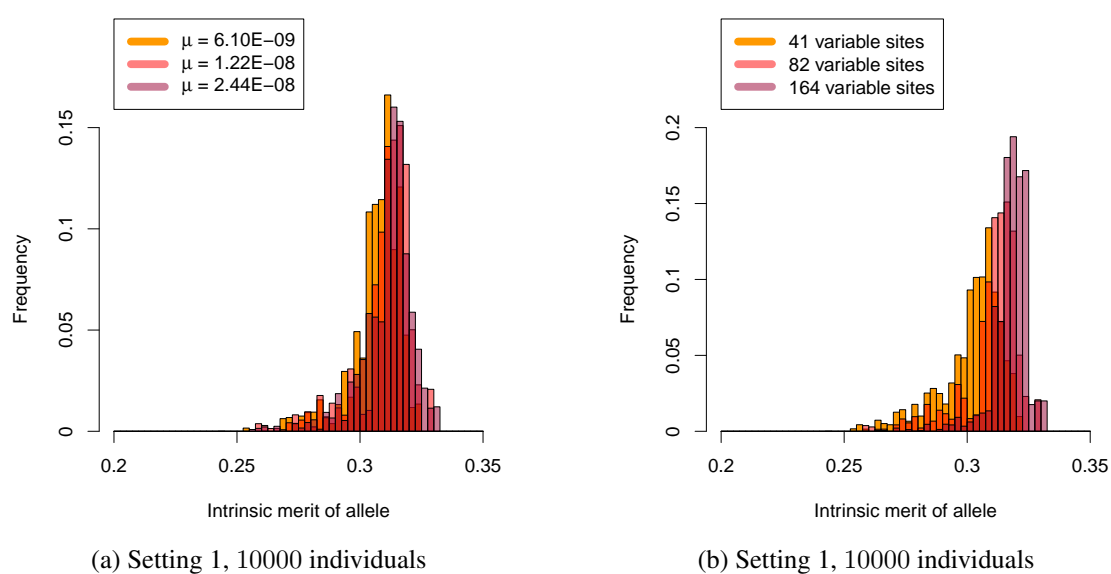
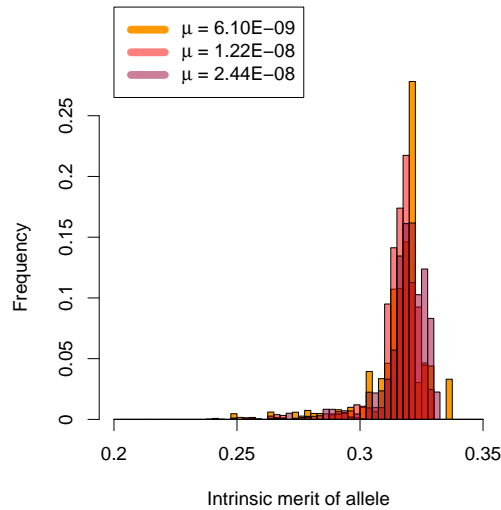
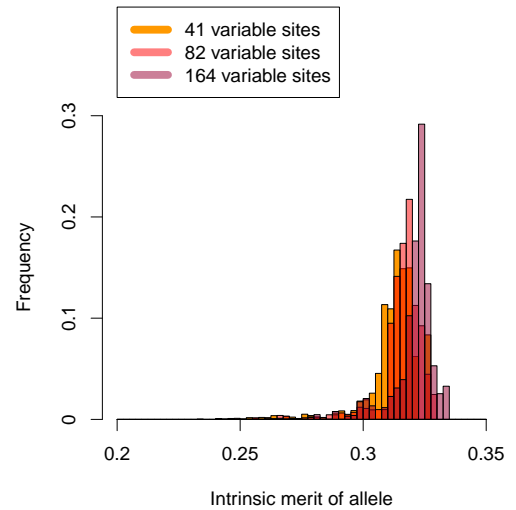


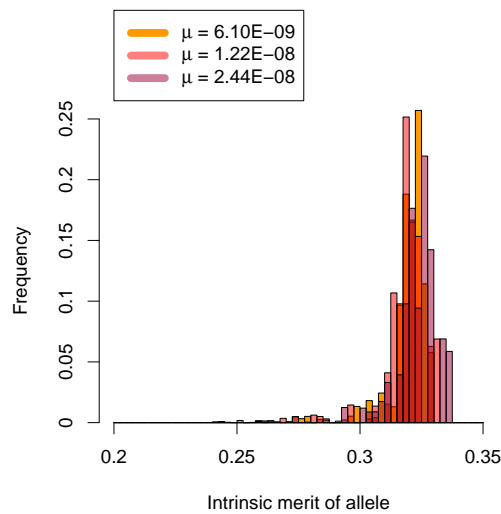
Figure B.63: Frequency distribution of allele intrinsic merits for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 10000.



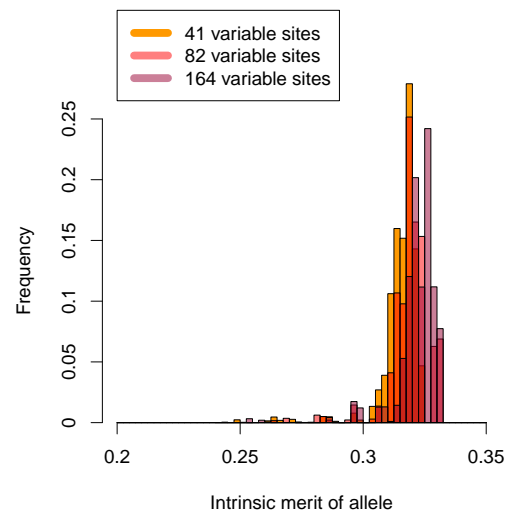
(a) Setting 1, 1 million individuals



(b) Setting 1, 1 million individuals



(c) Setting 1, 10 million individuals



(d) Setting 1, 10 million individuals

Figure B.64: Frequency distribution of allele intrinsic merits for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 1 million (top) and 10 millions (bottom).

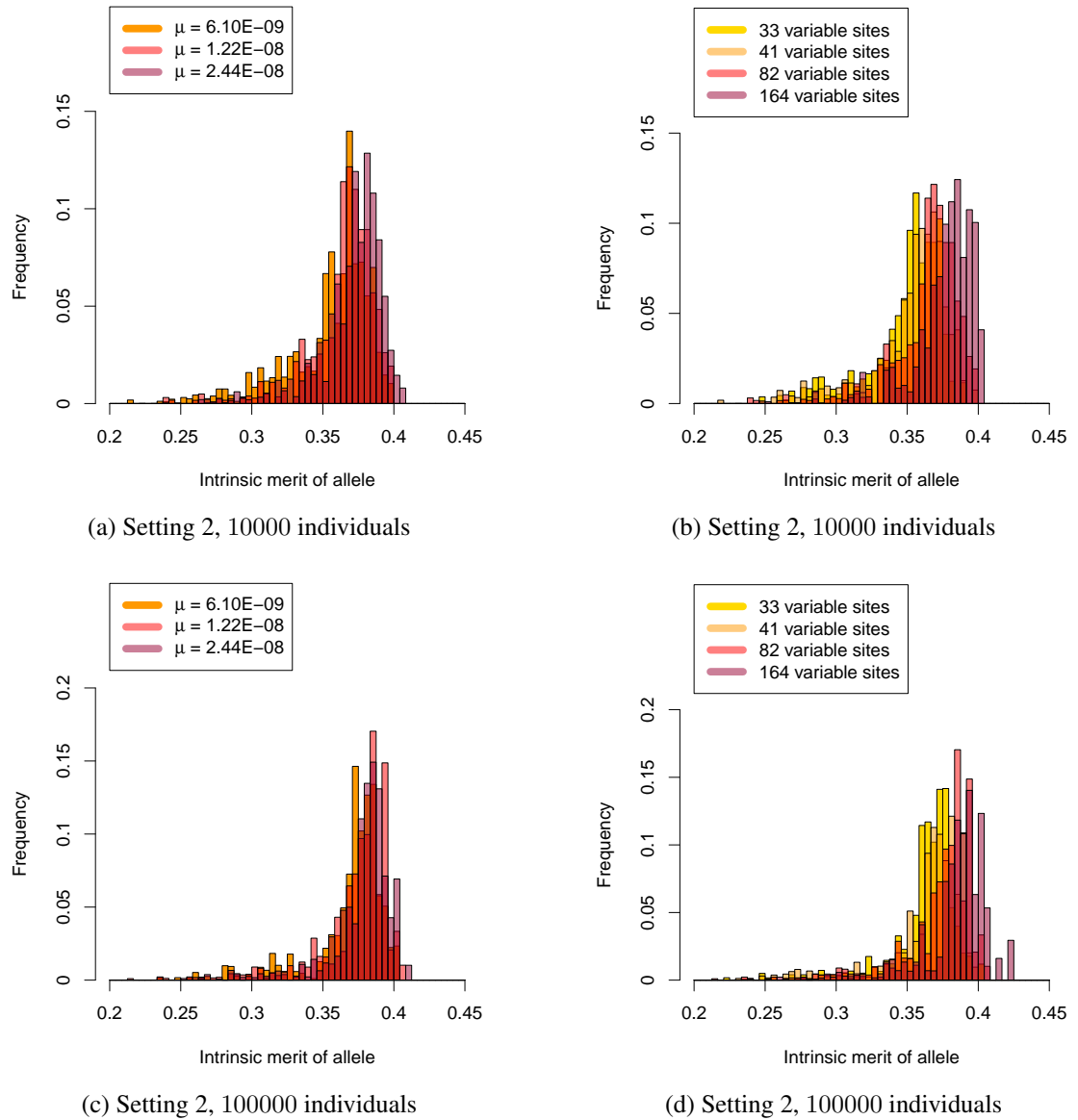
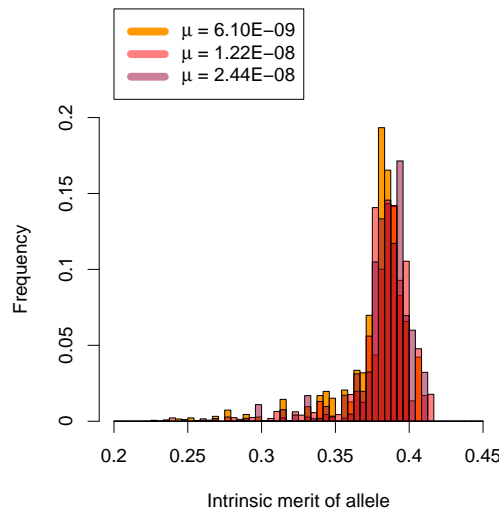
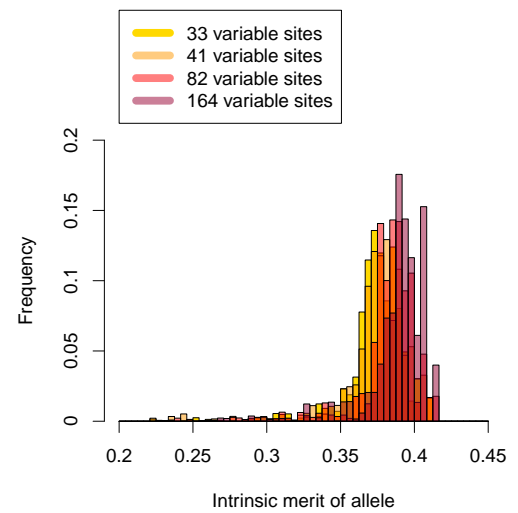


Figure B.65: Frequency distribution of allele intrinsic merits for setting 2 and different mutation rates (left) and number of variable sites (right) for a population size of 10000 (top) and 100000 (bottom).

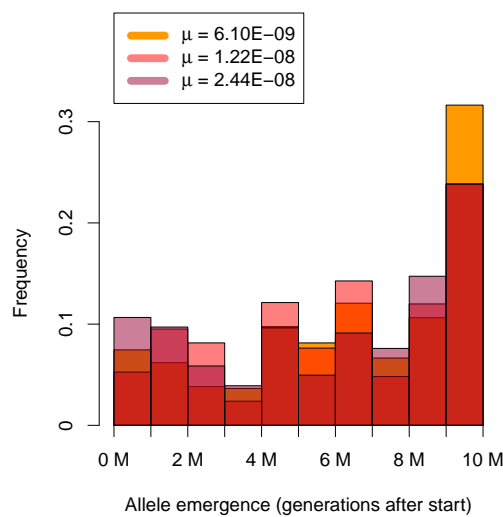


(a) Setting 2, 1 million individuals

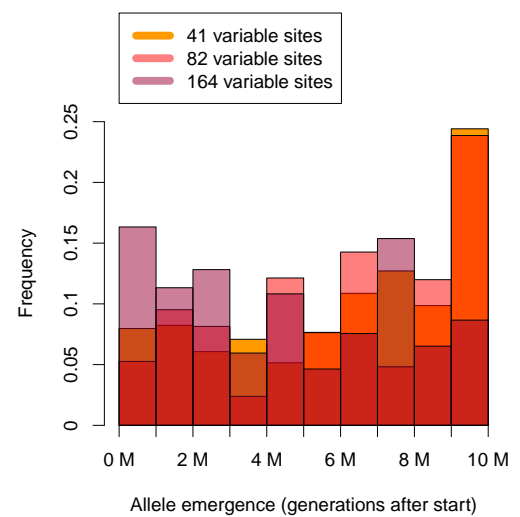


(b) Setting 2, 1 million individuals

Figure B.66: Frequency distribution of allele intrinsic merits for setting 2 and different mutation rates (left) and number of variable sites (right) for a population size of 1 million.

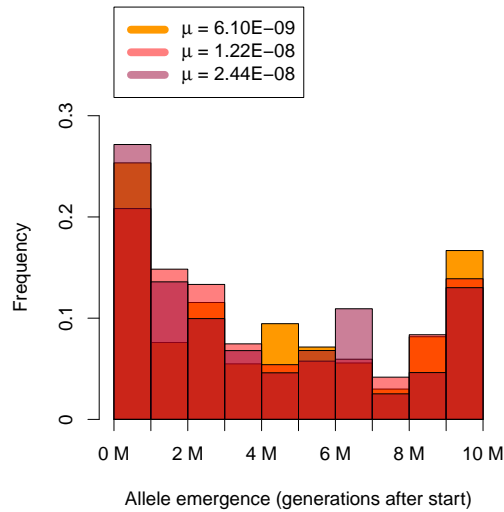


(a) Setting 1, 10000 individuals

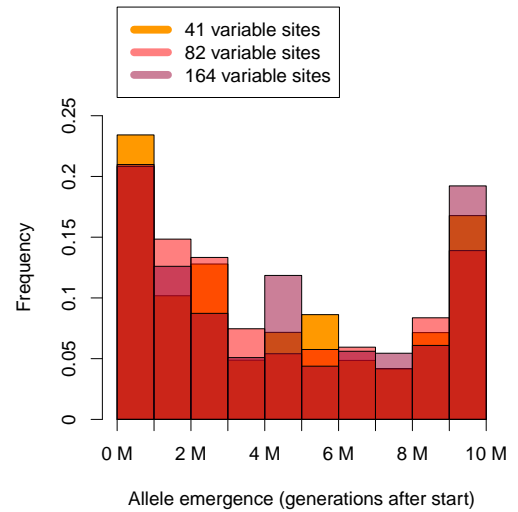


(b) Setting 1, 10000 individuals

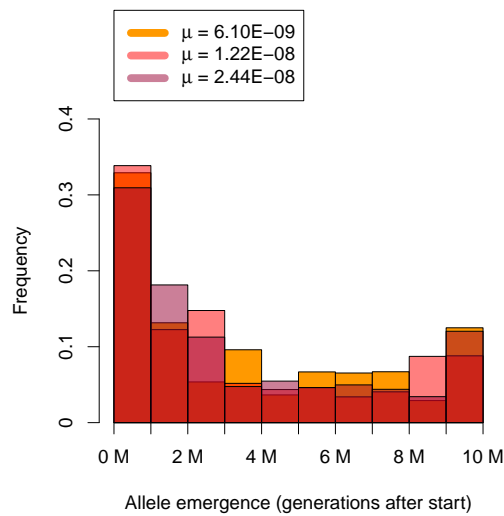
Figure B.67: Frequency distribution of emergence times of alleles for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 10000.



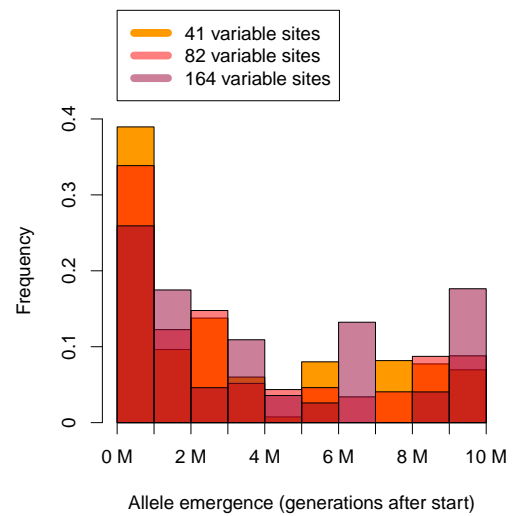
(a) Setting 1, 1 million individuals



(b) Setting 1, 1 million individuals

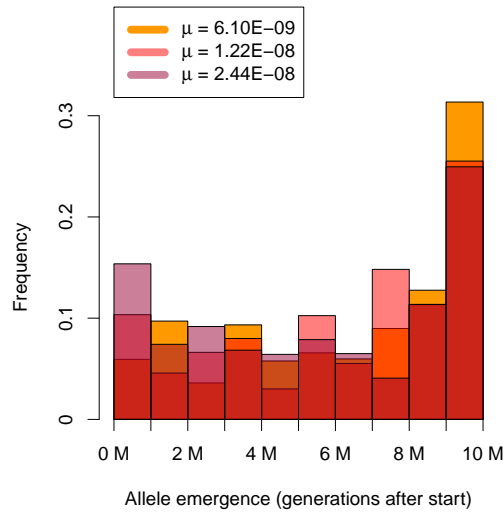


(c) Setting 1, 10 million individuals

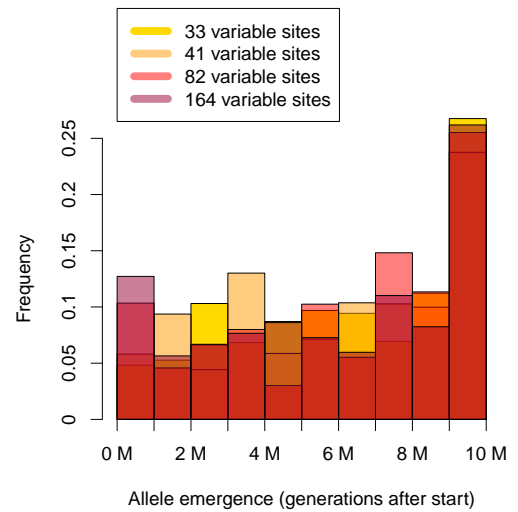


(d) Setting 1, 10 million individuals

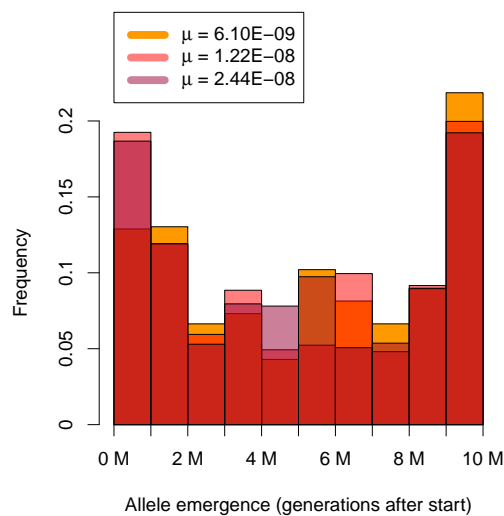
Figure B.68: Frequency distribution of emergence times of alleles for setting 1 and different mutation rates (left) and number of variable sites (right) for a population size of 1 million (top) and 10 millions (bottom).



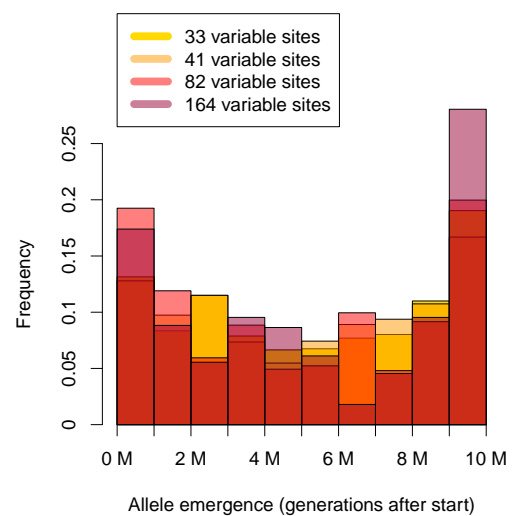
(a) Setting 2, 10000 individuals



(b) Setting 2, 10000 individuals



(c) Setting 2, 100000 individuals



(d) Setting 2, 100000 individuals

Figure B.69: Frequency distribution of emergence times of alleles for setting 2 and different mutation rates (left) and number of variable sites (right) for a population size of 10000 (top) and 100000 (bottom).

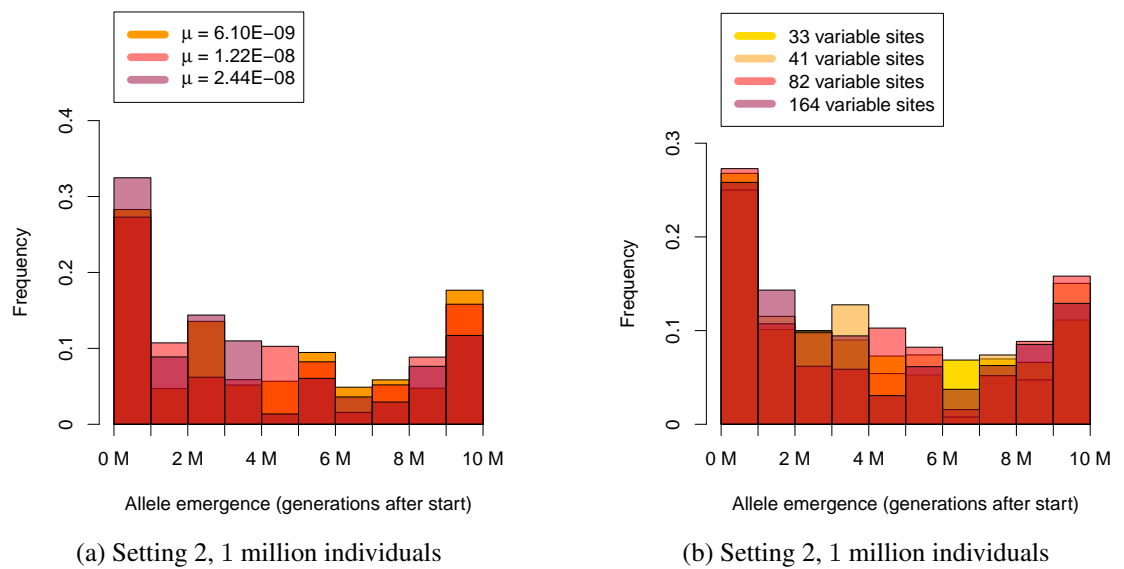
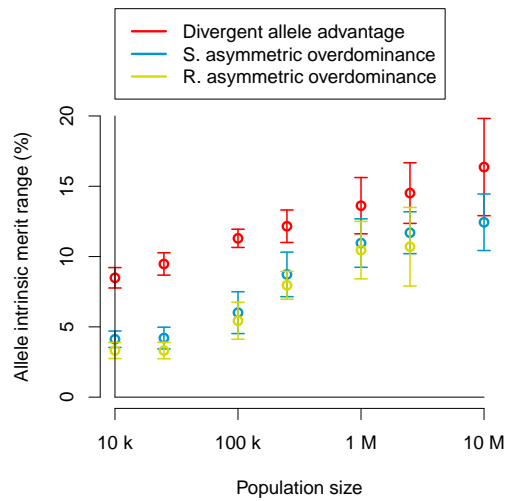


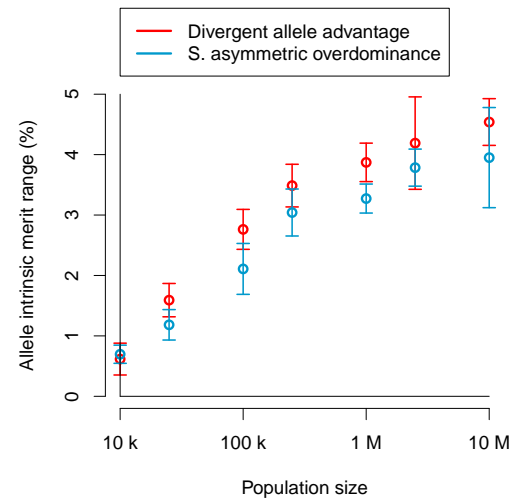
Figure B.70: Frequency distribution of emergence times of alleles for setting 2 and different mutation rates (left) and number of variable sites (right) for a population size of 1 million.

Appendix C

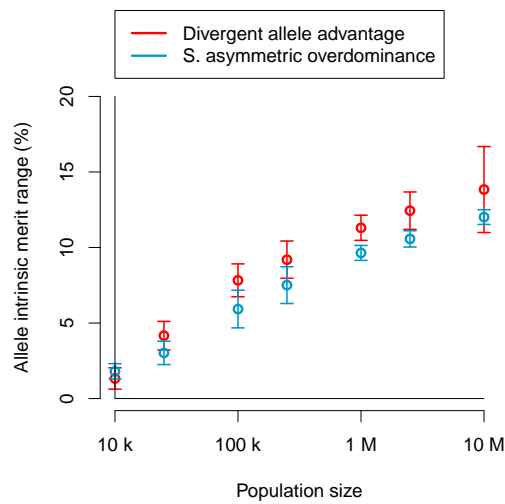
Additional figures for chapter 5



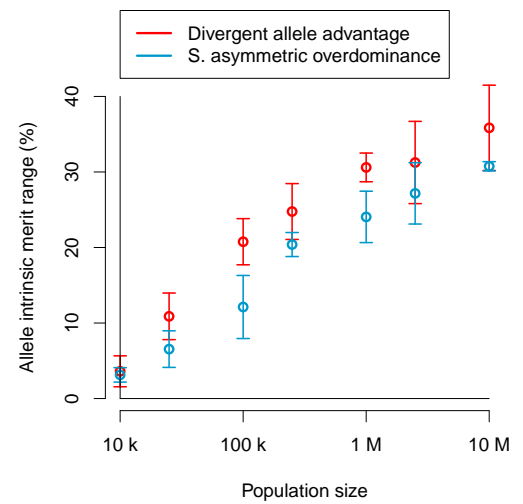
(a) Setting 3



(b) Setting 4

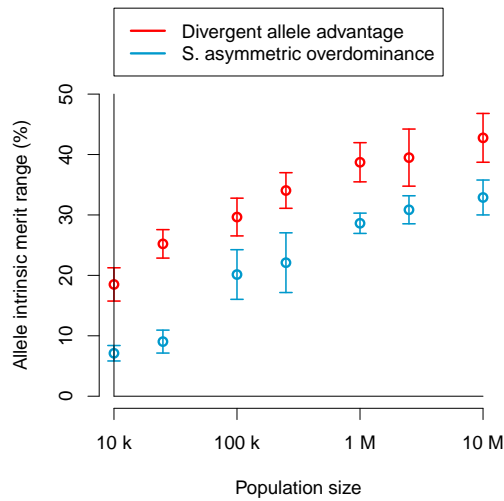


(c) Setting 5

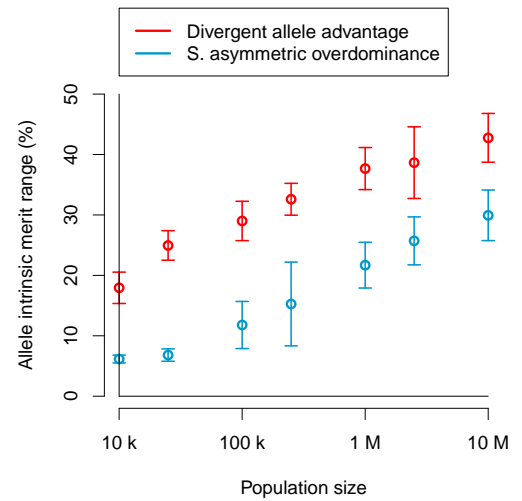


(d) Setting 6

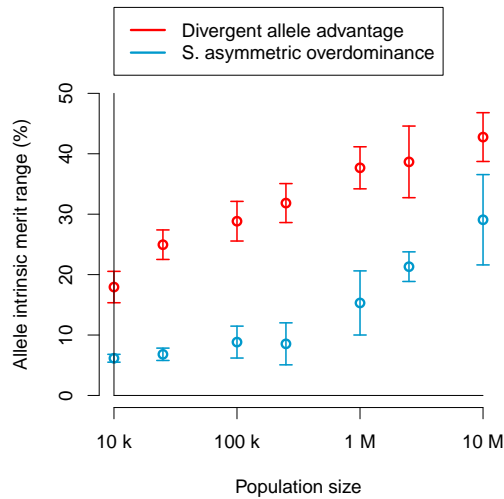
Figure C.1: Allele intrinsic merit range vs. population size (taking all alleles at the end of the simulation period into account) for settings 3 (top left), 4 (top right), 5 (bottom left) and 6 (bottom right).



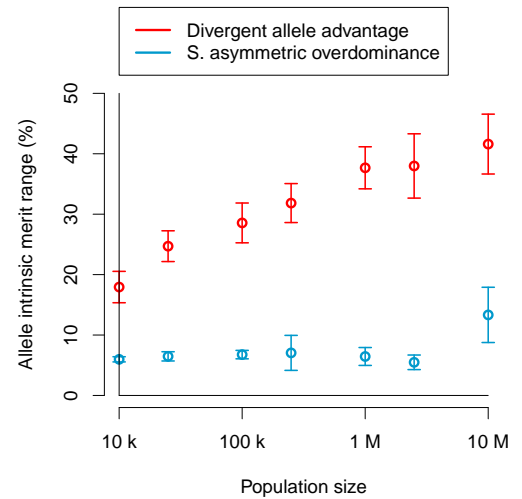
(a) No age threshold (all alleles)



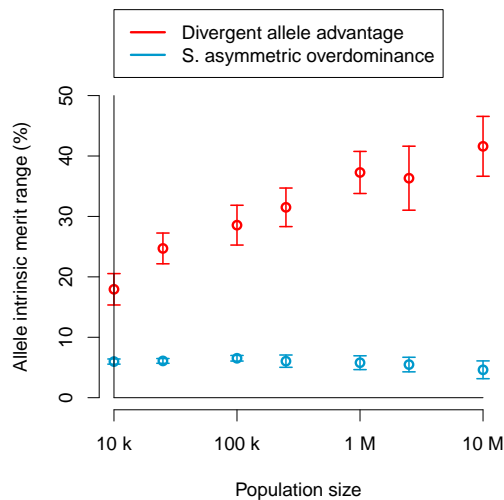
(b) Age threshold 10 generations



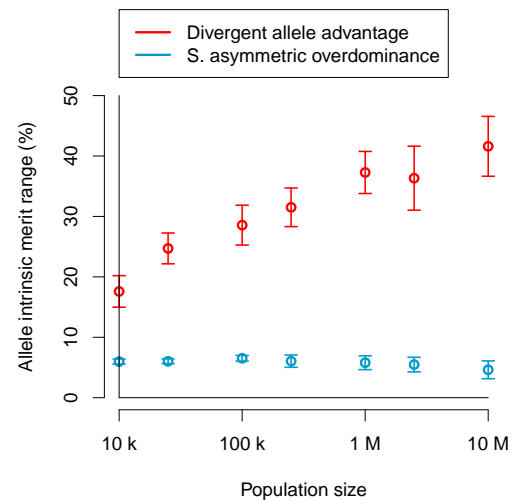
(c) Age threshold 25 generations



(d) Age threshold 100 generations

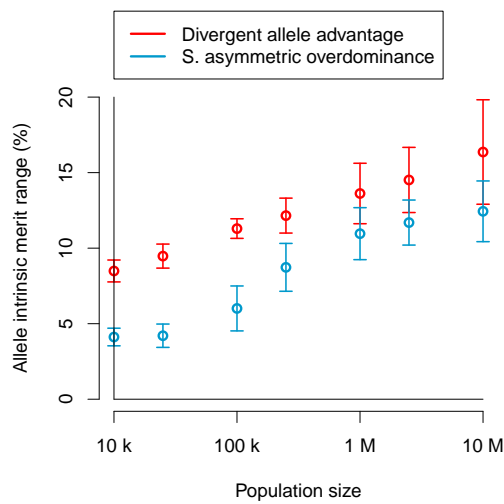


(e) Age threshold 250 generations

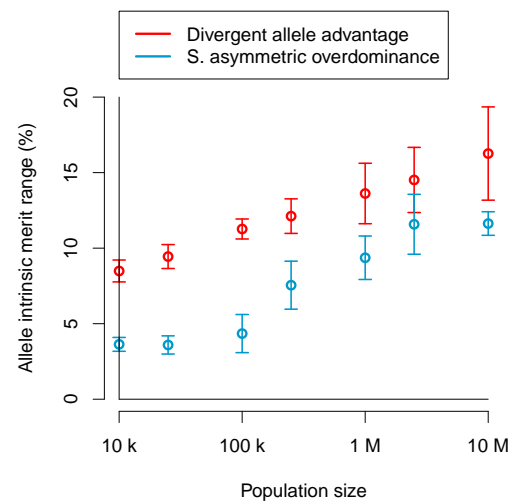


(f) Age threshold 1000 generations

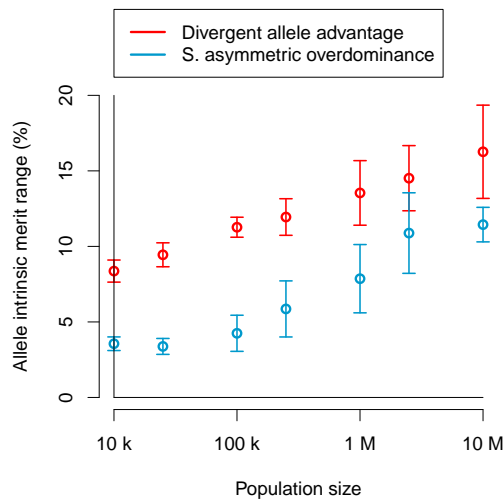
Figure C.2: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 2.



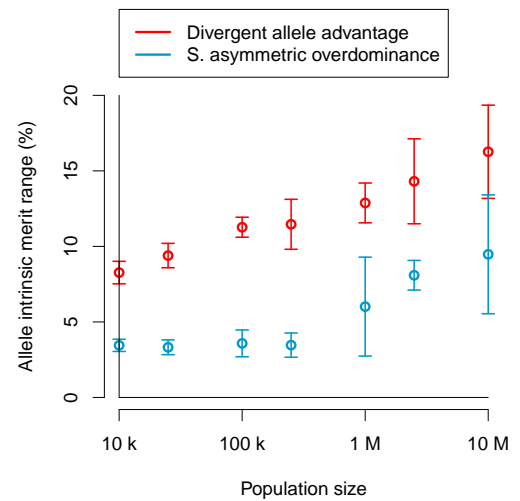
(a) No age threshold (all alleles)



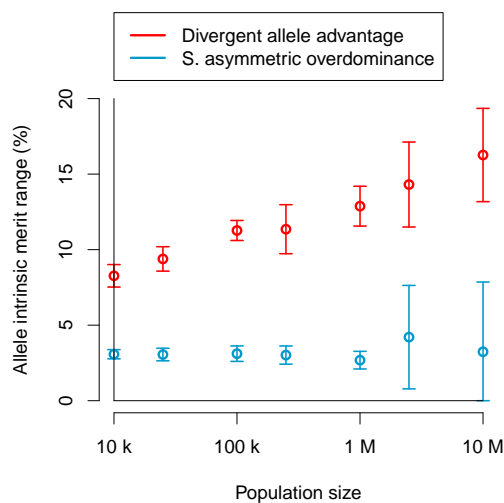
(b) Age threshold 10 generations



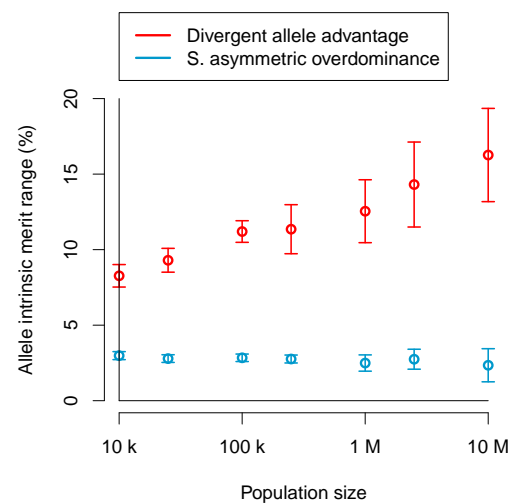
(c) Age threshold 25 generations



(d) Age threshold 100 generations

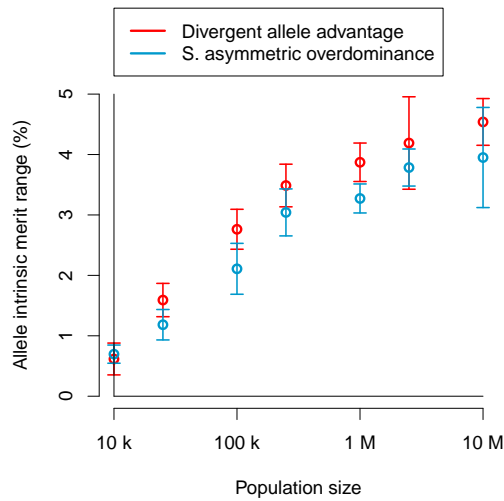


(e) Age threshold 250 generations

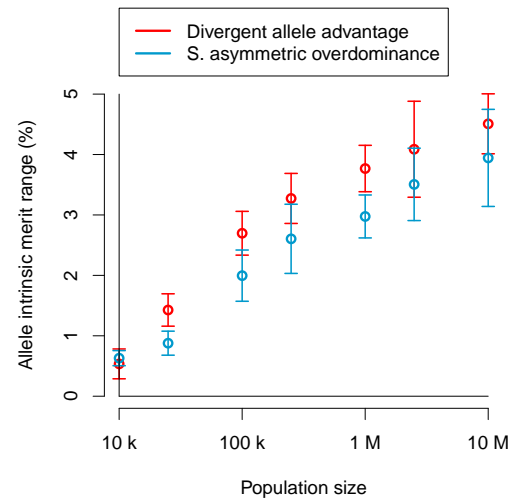


(f) Age threshold 1000 generations

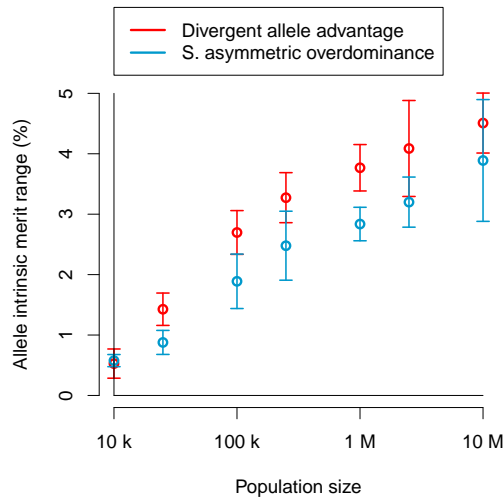
Figure C.3: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 3.



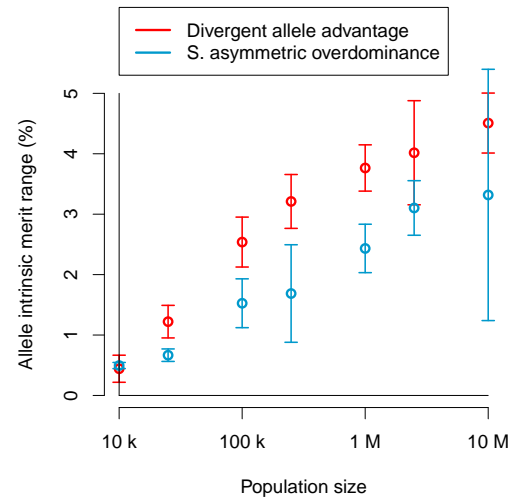
(a) No age threshold (all alleles)



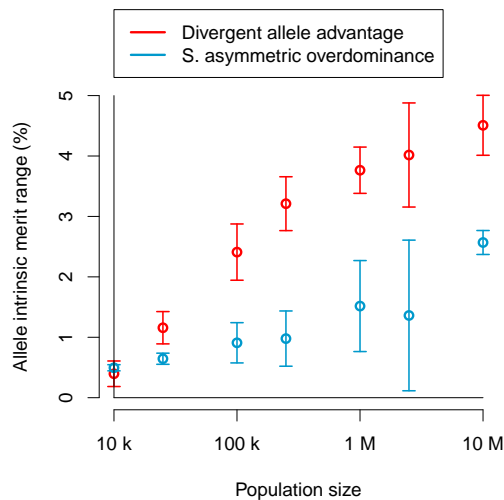
(b) Age threshold 10 generations



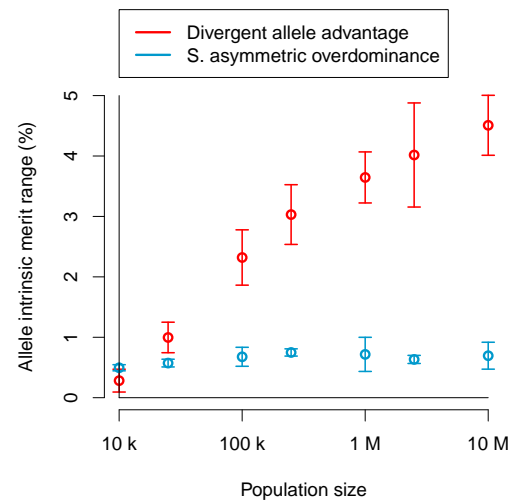
(c) Age threshold 25 generations



(d) Age threshold 100 generations

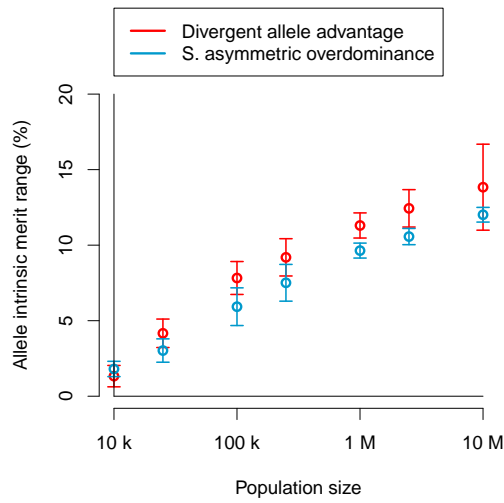


(e) Age threshold 250 generations

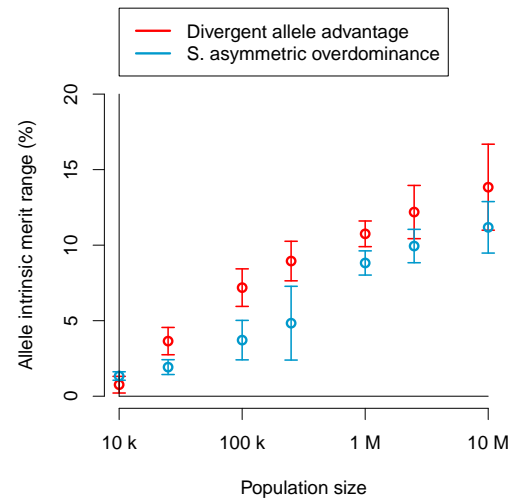


(f) Age threshold 1000 generations

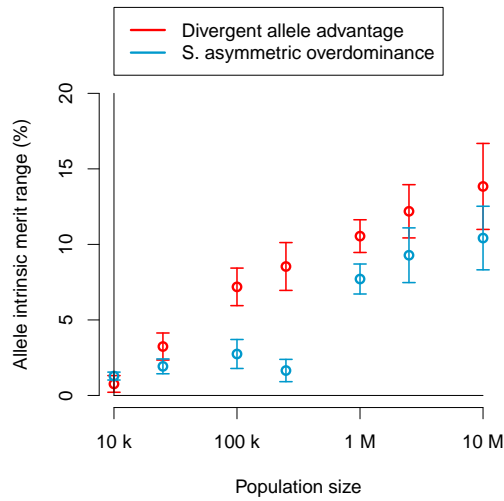
Figure C.4: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 4.



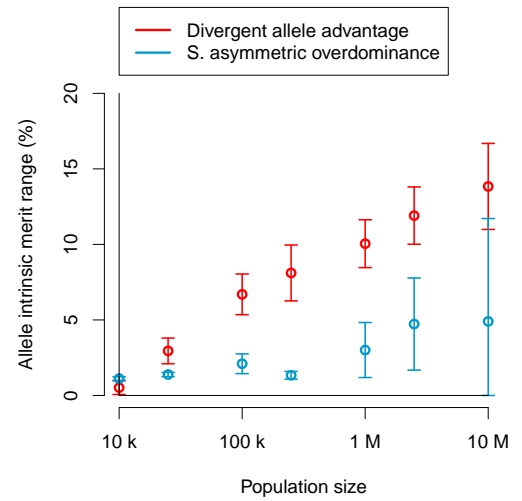
(a) No age threshold (all alleles)



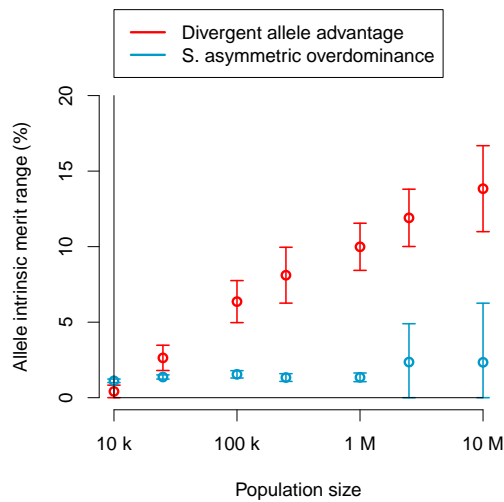
(b) Age threshold 10 generations



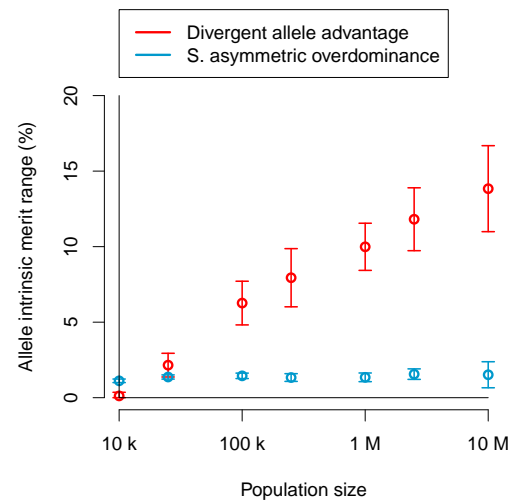
(c) Age threshold 25 generations



(d) Age threshold 100 generations

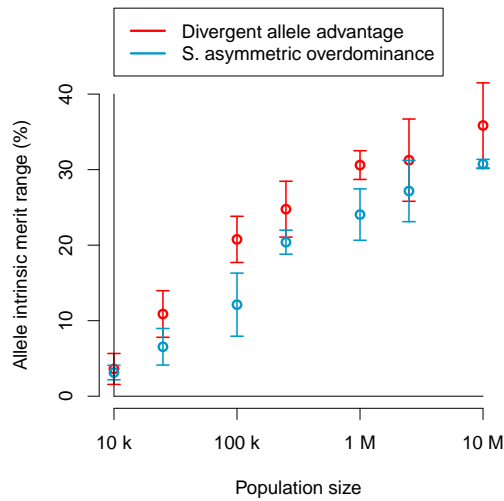


(e) Age threshold 250 generations

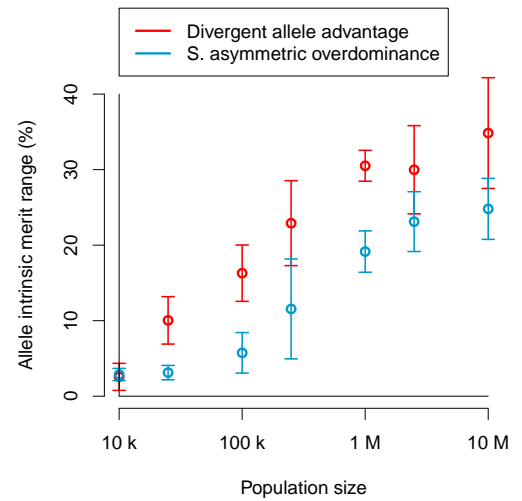


(f) Age threshold 1000 generations

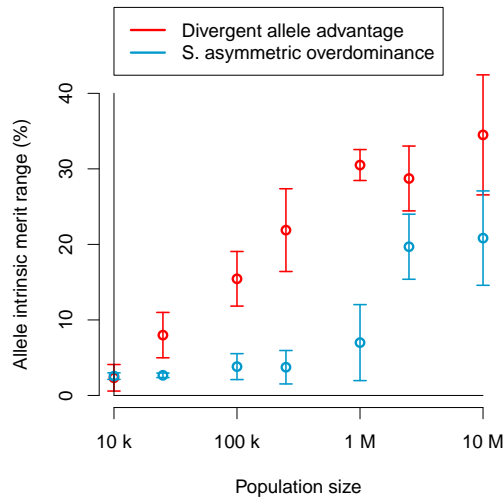
Figure C.5: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 5.



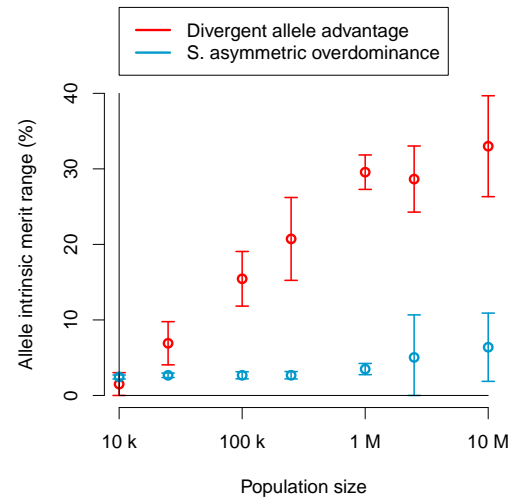
(a) No age threshold (all alleles)



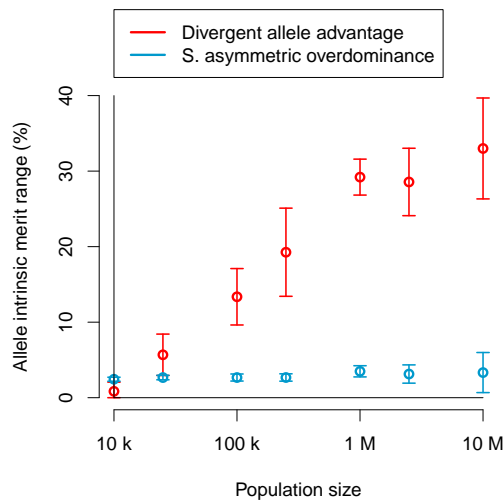
(b) Age threshold 10 generations



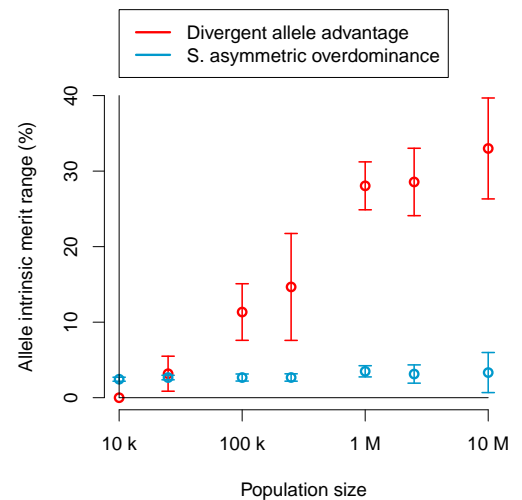
(c) Age threshold 25 generations



(d) Age threshold 100 generations

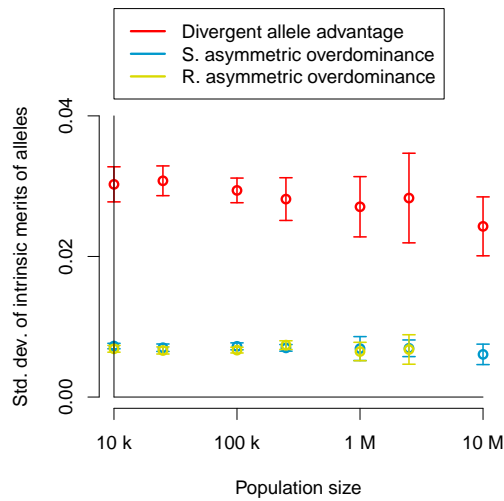


(e) Age threshold 250 generations

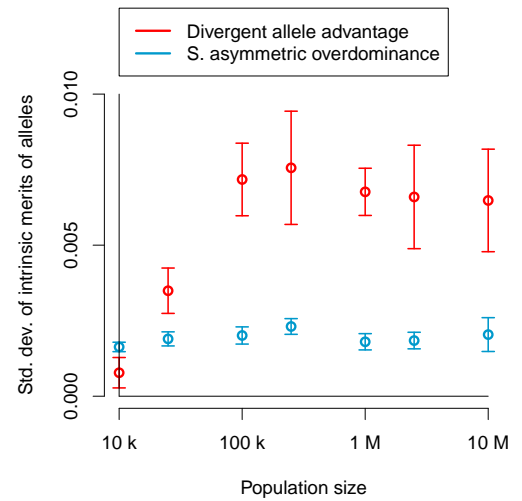


(f) Age threshold 1000 generations

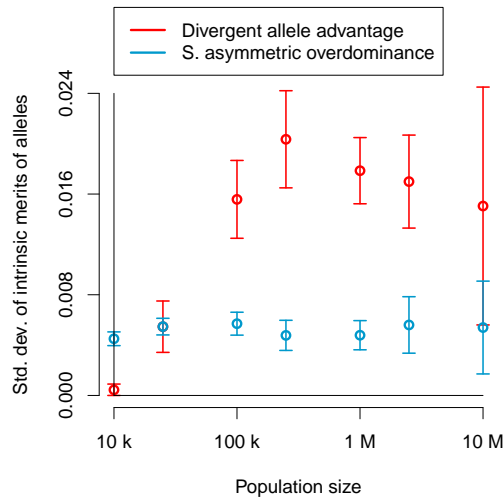
Figure C.6: Range of allele intrinsic merits when using different age thresholds to remove transient alleles. Results correspond to setting 6.



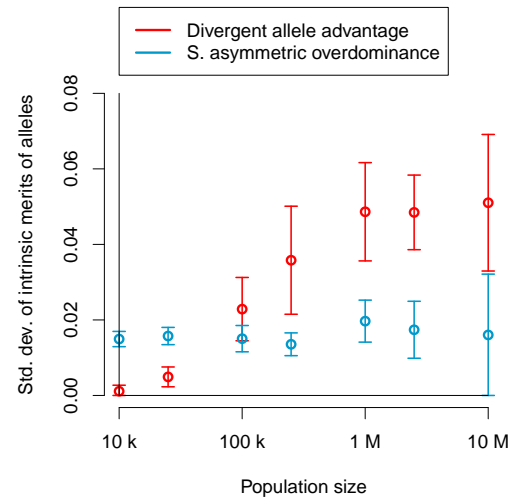
(a) 250 generations, setting 3



(b) 250 generations, setting 4

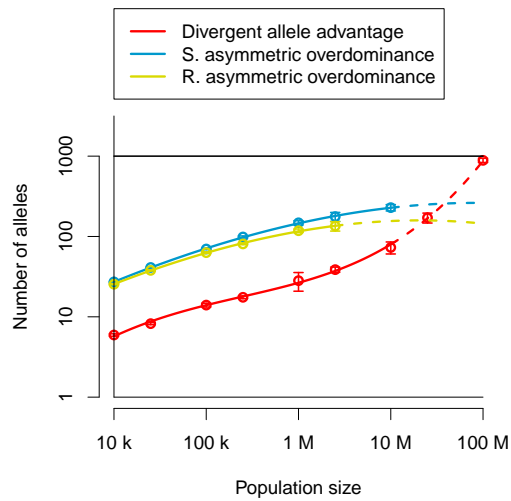


(c) 250 generations, setting 5

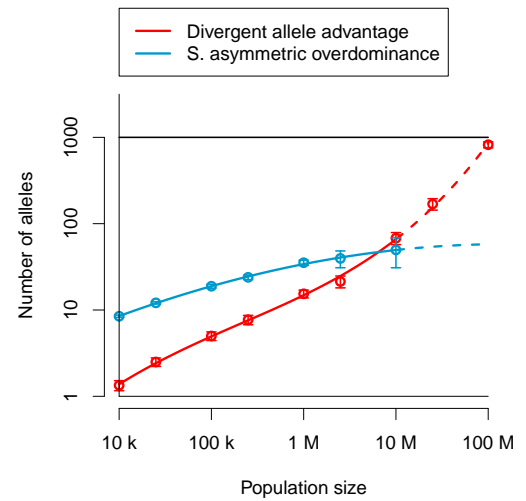


(d) 250 generations, setting 6

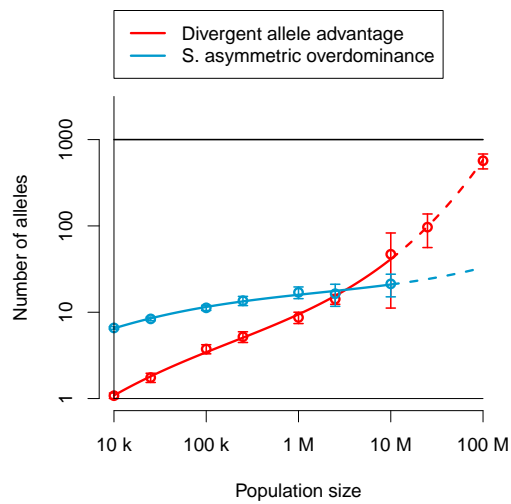
Figure C.7: Standard deviation of the intrinsic merits of the final set of alleles when using an age threshold of 250 generations to remove transient alleles. Results correspond to settings 3 (top left), 4 (top right), 5 (bottom left) and 6 (bottom right).



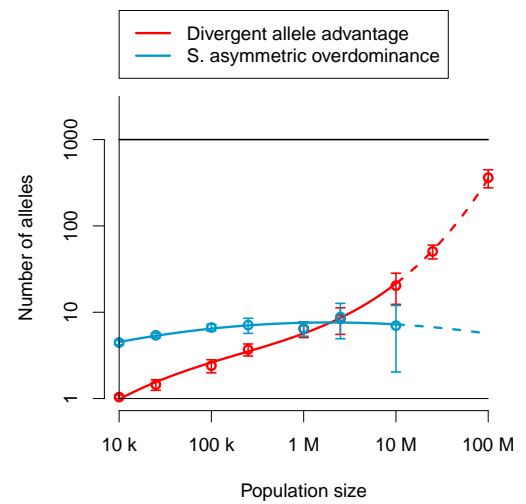
(a) 250 generations, setting 3



(b) 250 generations, setting 4

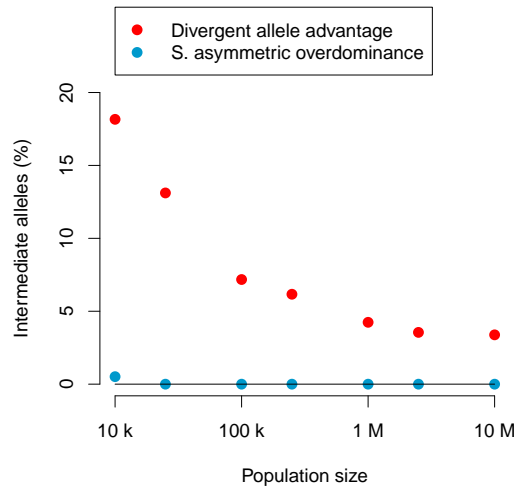


(c) 250 generations, setting 5

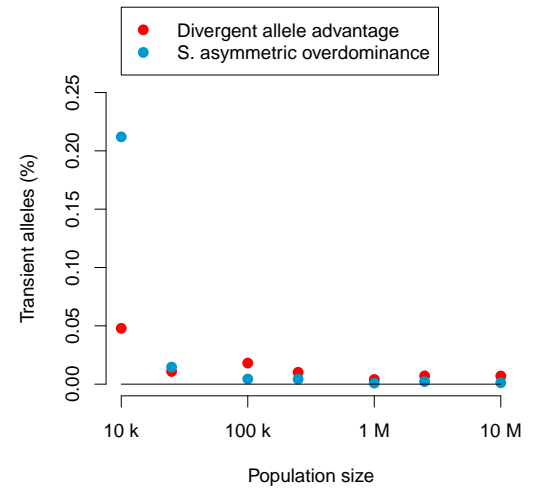


(d) 250 generations, setting 6

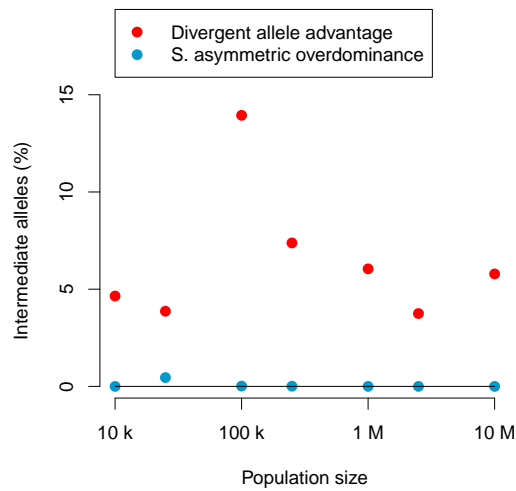
Figure C.8: Number of alleles at the end of the simulation when using an age threshold of 250 generations to remove transient alleles. Results correspond to settings 3 (top left), 4 (top right), 5 (bottom left) and 6 (bottom right).



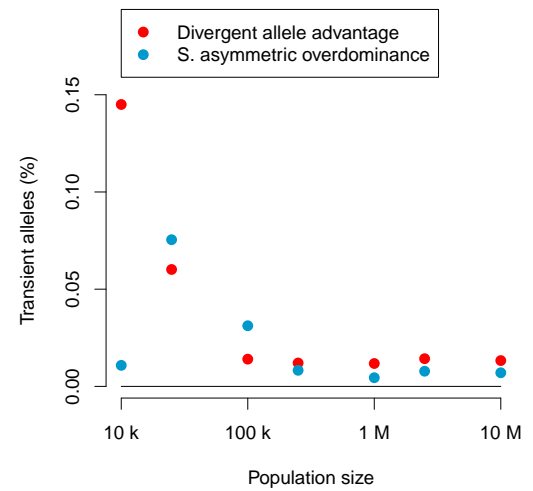
(a) Share of intermediate alleles, setting 3



(b) Share of transient alleles, setting 3

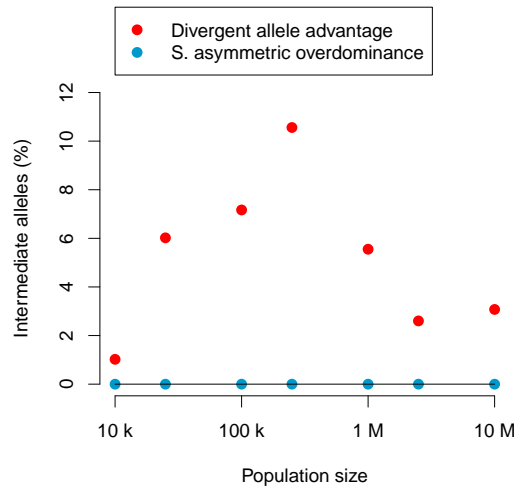


(c) Share of intermediate alleles, setting 4

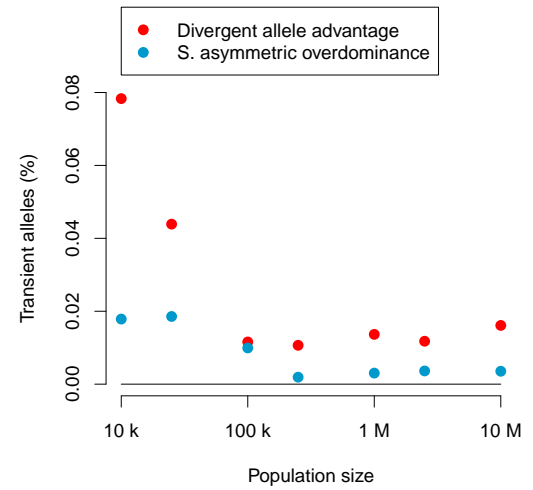


(d) Share of transient alleles, setting 4

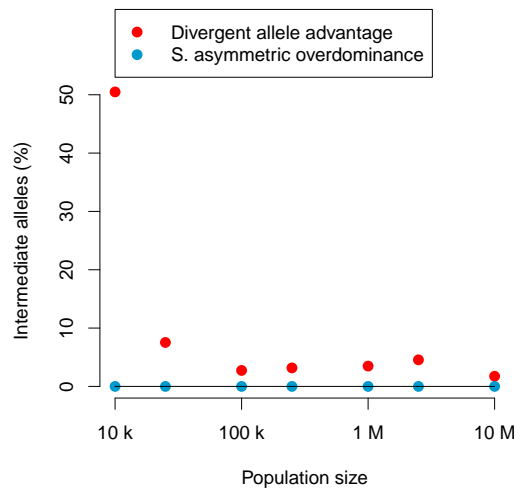
Figure C.9: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 3 (top) and 4 (bottom). The threshold for allele intrinsic merit was 0.265 for setting 3 and 0.256 for setting 4, and alleles that emerged in the last 250 generations were classified as transient.



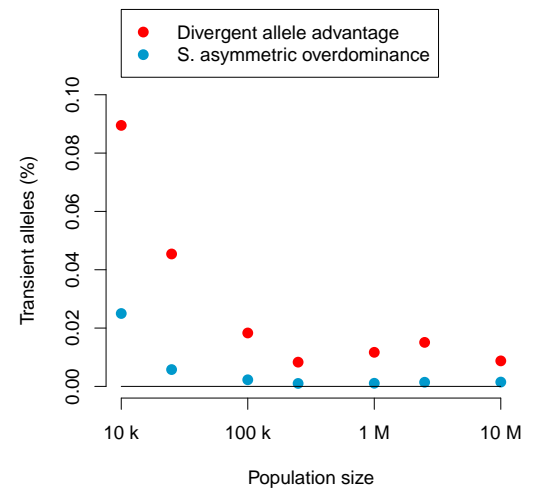
(a) Share of intermediate alleles, setting 5



(b) Share of transient alleles, setting 5

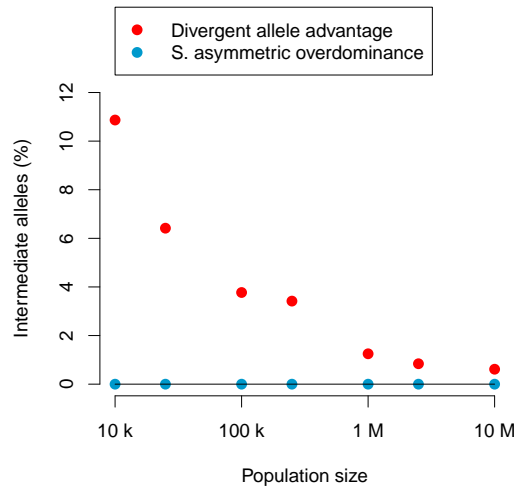


(c) Share of intermediate alleles, setting 6

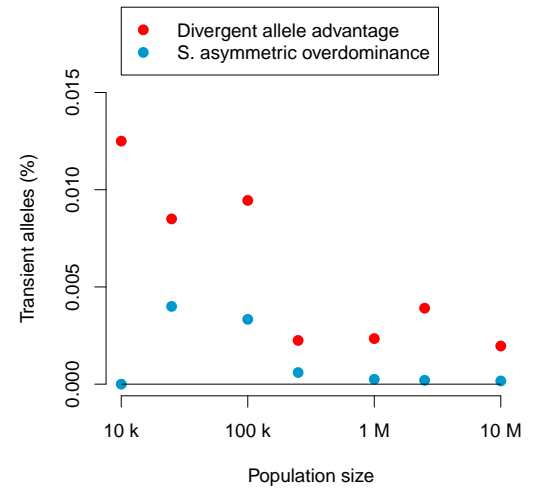


(d) Share of transient alleles, setting 6

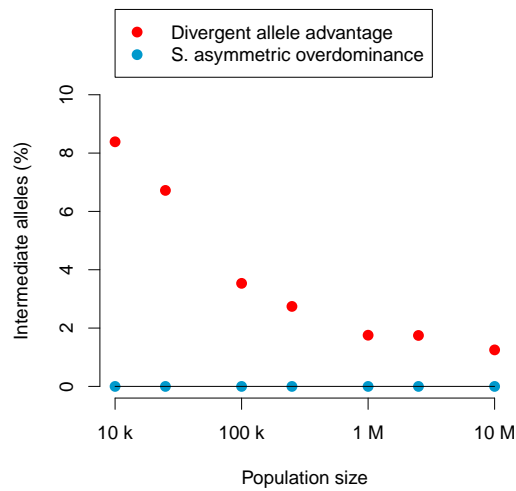
Figure C.10: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 5 (top) and 6 (bottom). The threshold for allele intrinsic merit was 0.275 for setting 5 and 0.365 for setting 6, and alleles that emerged in the last 250 generations were classified as transient.



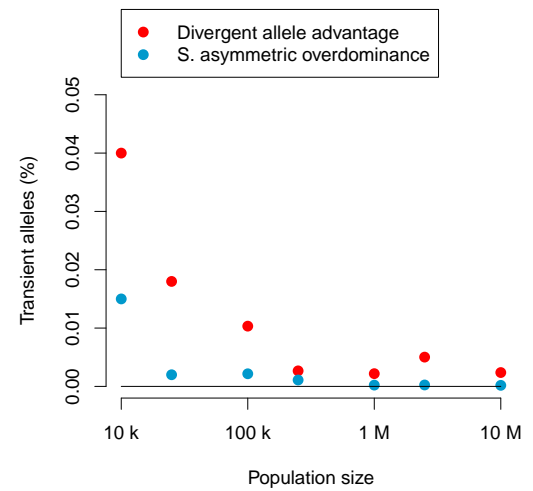
(a) Share of intermediate alleles, setting 1



(b) Share of transient alleles, setting 1

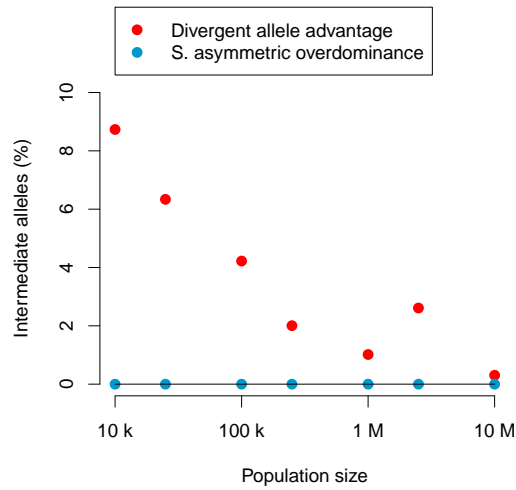


(c) Share of intermediate alleles, setting 2

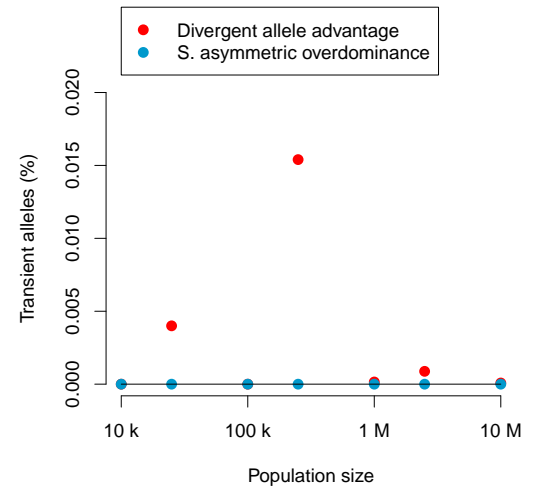


(d) Share of transient alleles, setting 2

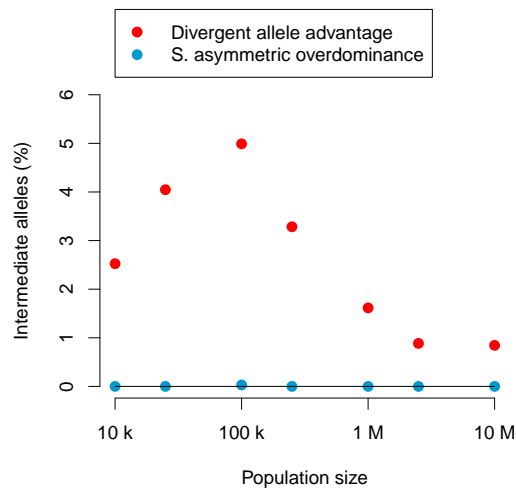
Figure C.11: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 1 (top) and 2 (bottom). The threshold for allele intrinsic merit was 0.275 for setting 1 and 0.3 for setting 2, and alleles that emerged in the last 250 generations were classified as transient.



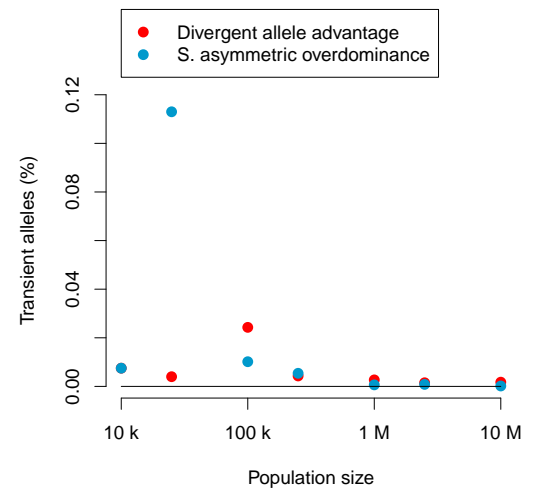
(a) Share of intermediate alleles, setting 3



(b) Share of transient alleles, setting 3

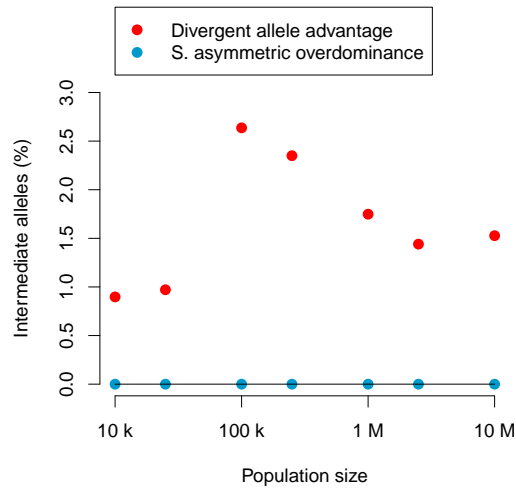


(c) Share of intermediate alleles, setting 4

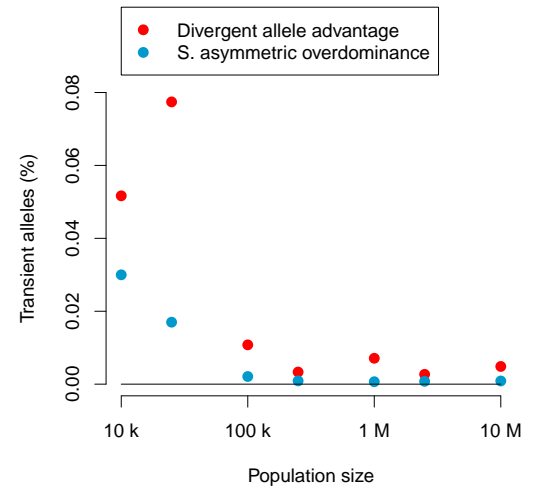


(d) Share of transient alleles, setting 4

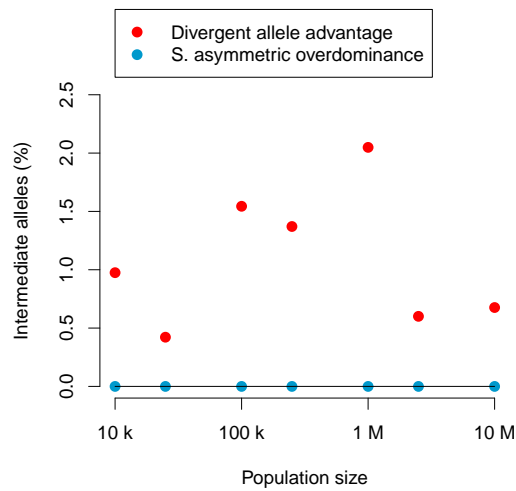
Figure C.12: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 3 (top) and 4 (bottom). The threshold for allele intrinsic merit was 0.25 for setting 3 and 0.253 for setting 4, and alleles that emerged in the last 250 generations were classified as transient.



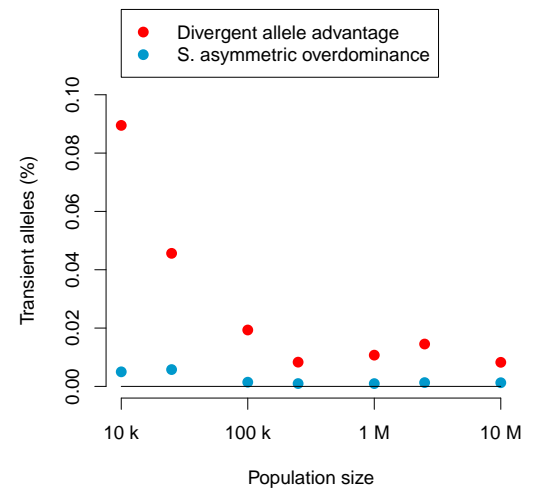
(a) Share of intermediate alleles, setting 5



(b) Share of transient alleles, setting 5

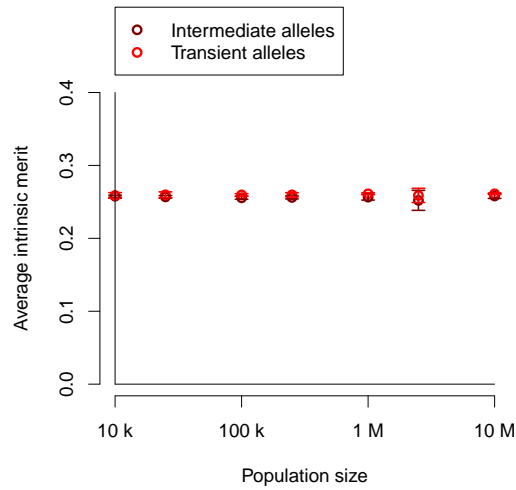


(c) Share of intermediate alleles, setting 6

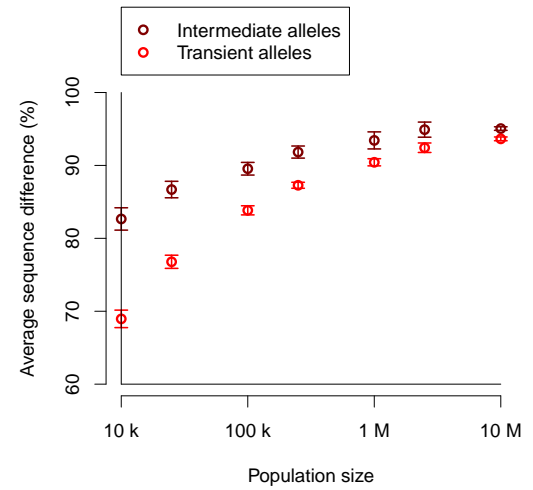


(d) Share of transient alleles, setting 6

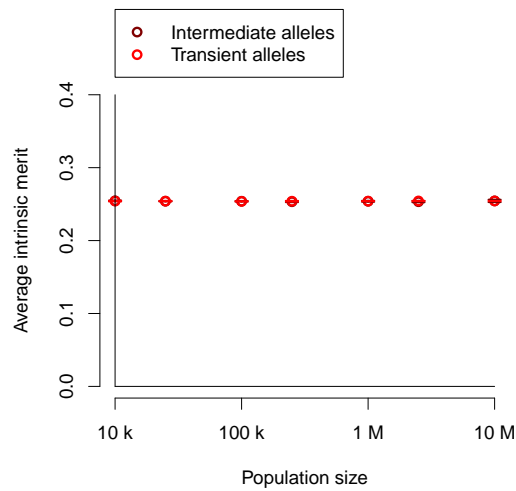
Figure C.13: Share of the gene pool of intermediate and transient alleles in the DAA and SAsOD models for settings 5 (top) and 6 (bottom). The threshold for allele intrinsic merit was 0.265 for setting 5 and 0.35 for setting 6, and alleles that emerged in the last 250 generations were classified as transient.



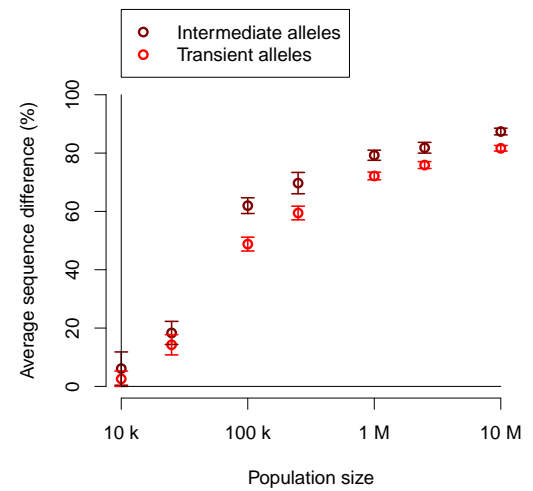
(a) Setting 3, less restrictive threshold



(b) Setting 3, less restrictive threshold



(c) Setting 4, less restrictive threshold



(d) Setting 4, less restrictive threshold

Figure C.14: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 3 (top) and 4 (bottom). The threshold for allele intrinsic merit was 0.265 for setting 3 and 0.256 for setting 4, and alleles emerged in the last 250 generations were classified as transient.

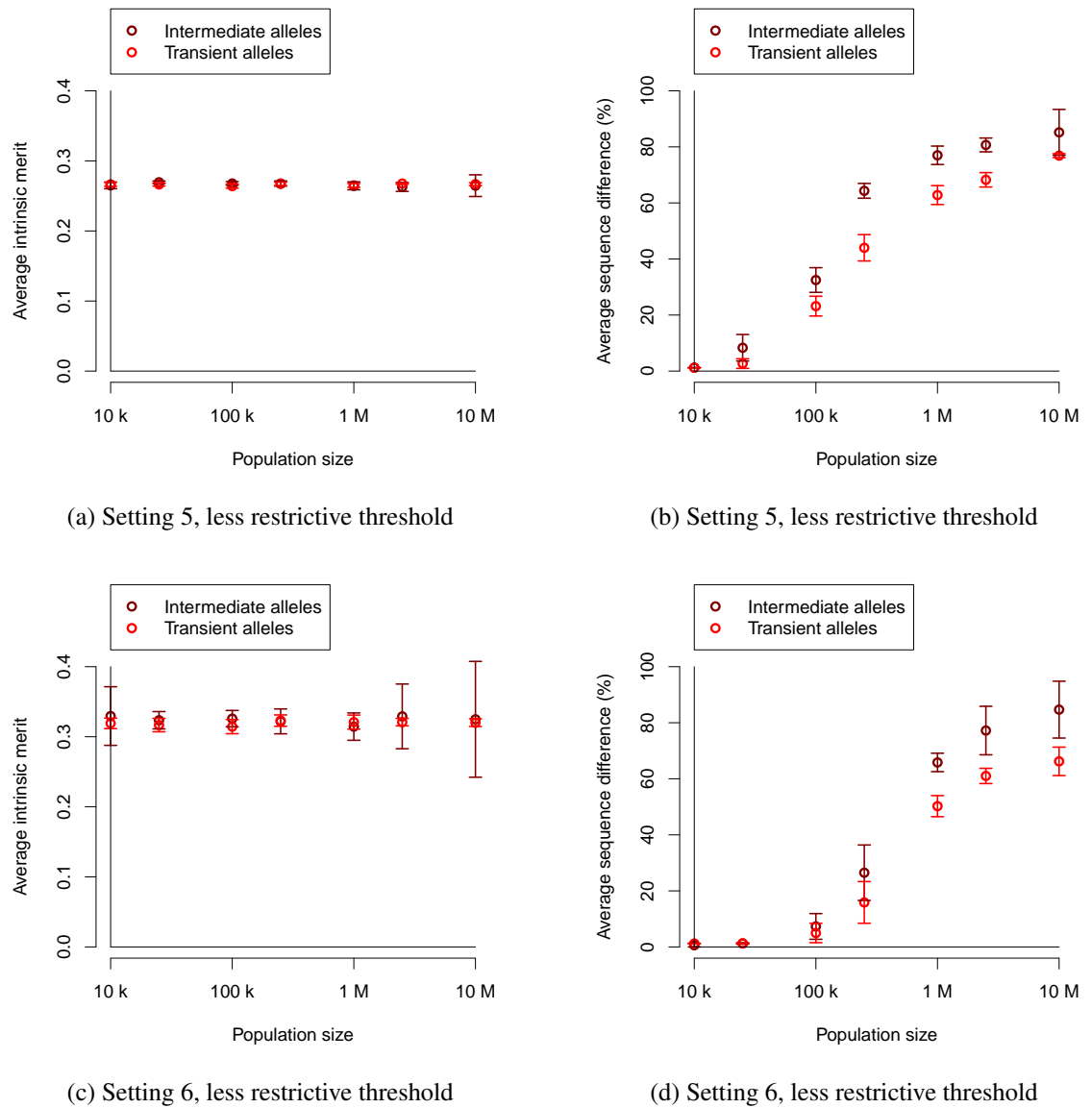
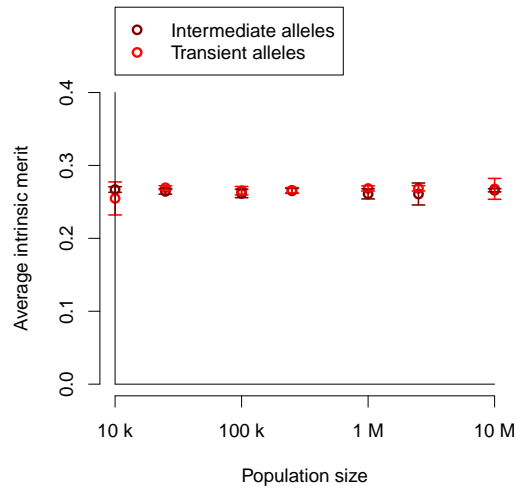
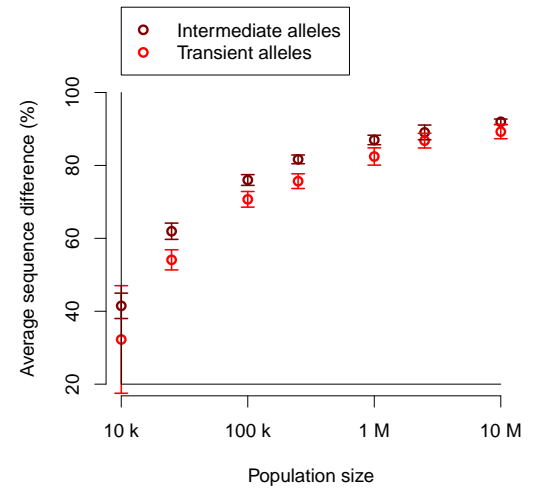


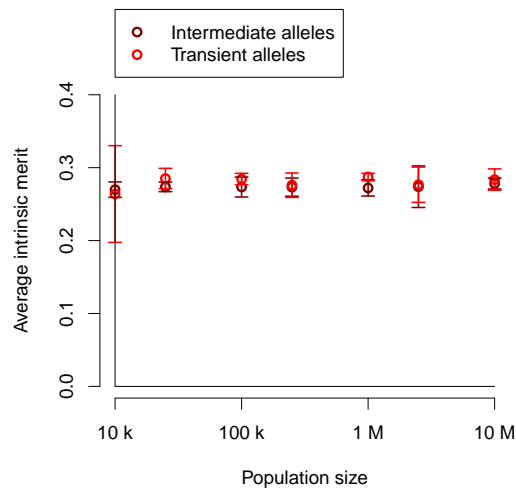
Figure C.15: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 5 (top) and 6 (bottom). The threshold for allele intrinsic merit was 0.275 for setting 5 and 0.365 for setting 6, and alleles emerged in the last 250 generations were classified as transient.



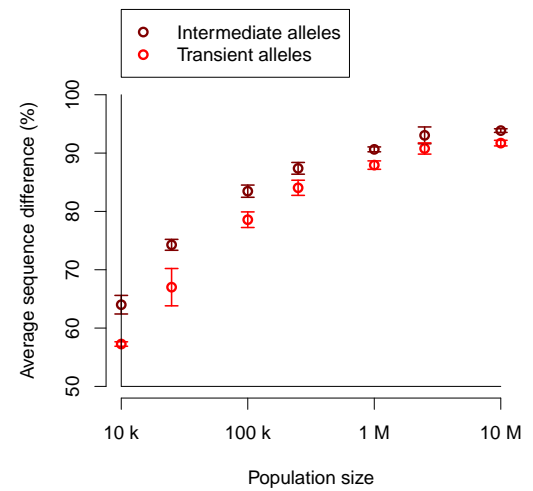
(a) Setting 1, stricter threshold



(b) Setting 1, stricter threshold

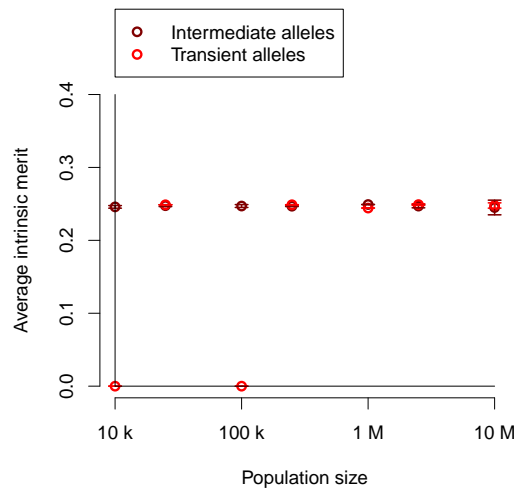


(c) Setting 2, stricter threshold

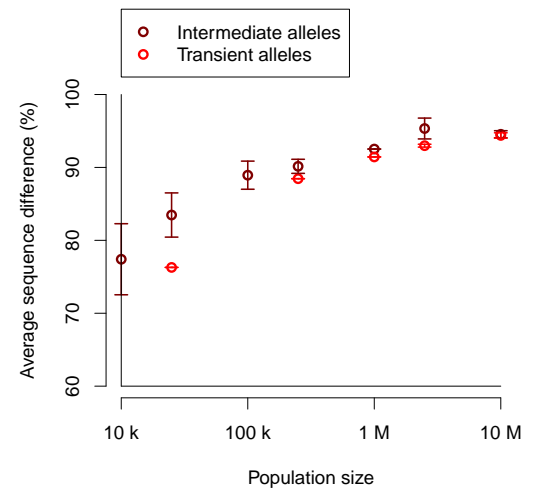


(d) Setting 2, stricter threshold

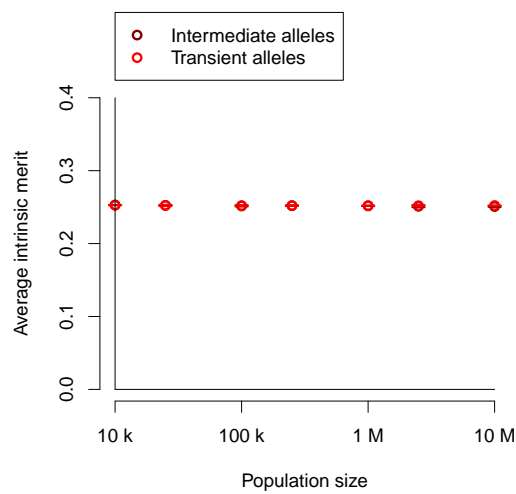
Figure C.16: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 1 (top) and 2 (bottom). The threshold for allele intrinsic merit was 0.275 for setting 1 and 0.3 for setting 2, and alleles emerged in the last 250 generations were classified as transient.



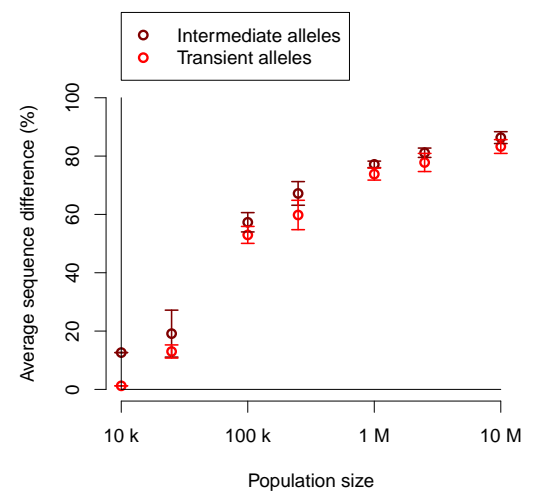
(a) Setting 3, stricter threshold



(b) Setting 3, stricter threshold

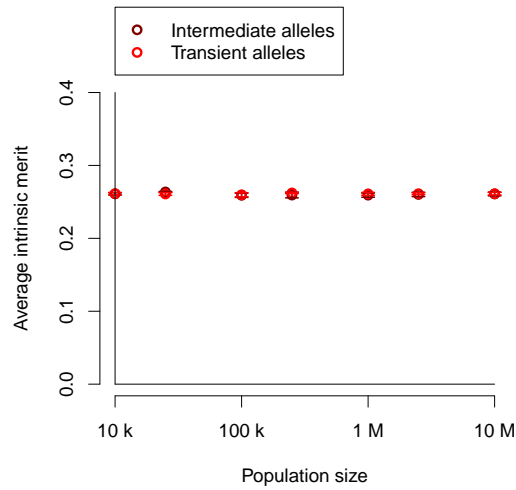


(c) Setting 4, stricter threshold

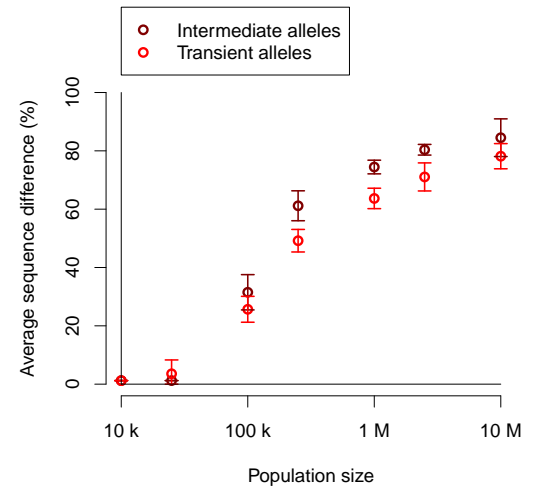


(d) Setting 4, stricter threshold

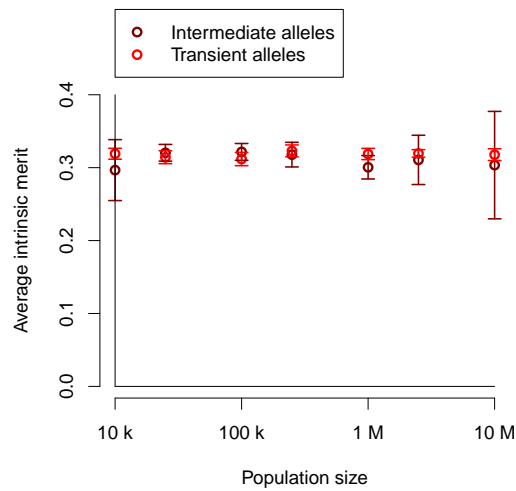
Figure C.17: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 3 (top) and 4 (bottom). The threshold for allele intrinsic merit was 0.25 for setting 3 and 0.253 for setting 4, and alleles emerged in the last 250 generations were classified as transient.



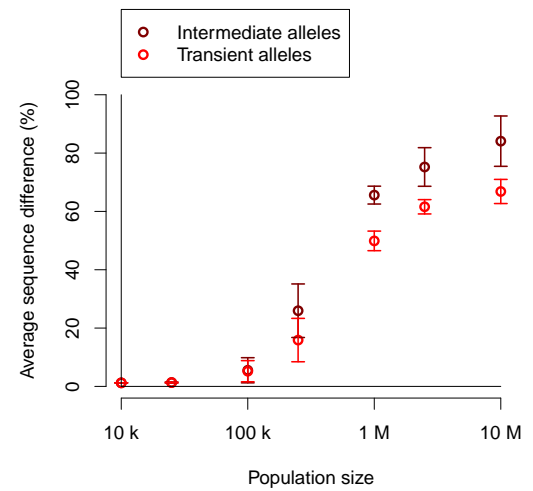
(a) Setting 5, stricter threshold



(b) Setting 5, stricter threshold



(c) Setting 6, stricter threshold



(d) Setting 6, stricter threshold

Figure C.18: Intrinsic merit and amino acid sequence difference for intermediate and transient alleles in the DAA model for settings 5 (top) and 6 (bottom). The threshold for allele intrinsic merit was 0.265 for setting 5 and 0.350 for setting 6, and alleles emerged in the last 250 generations were classified as transient.

Appendix D

Additional figures for chapter 6

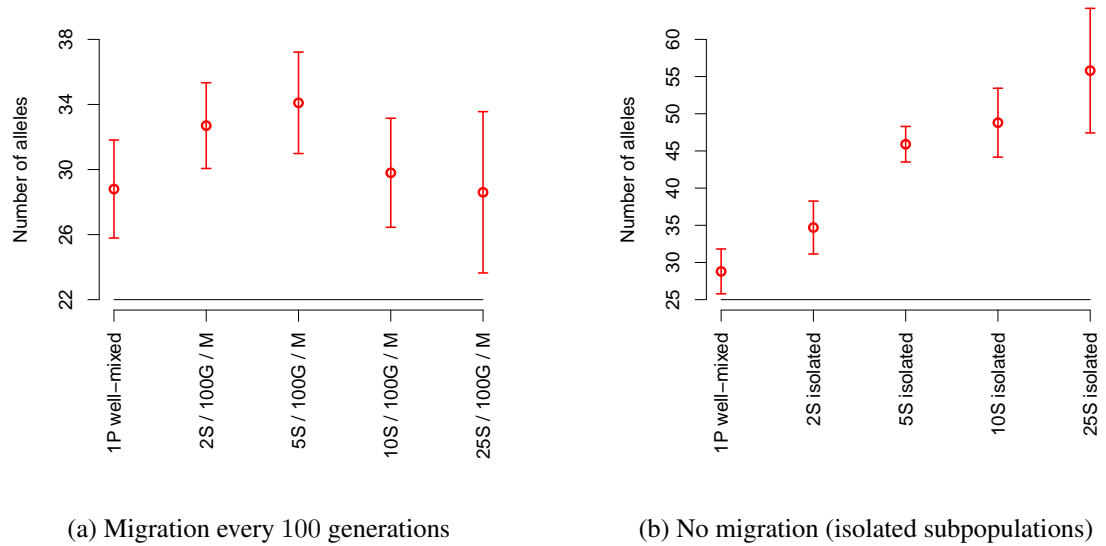
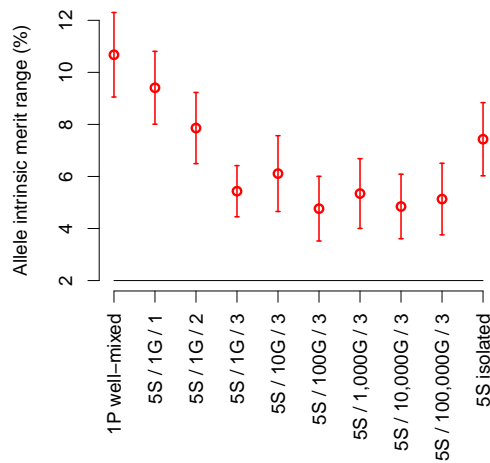
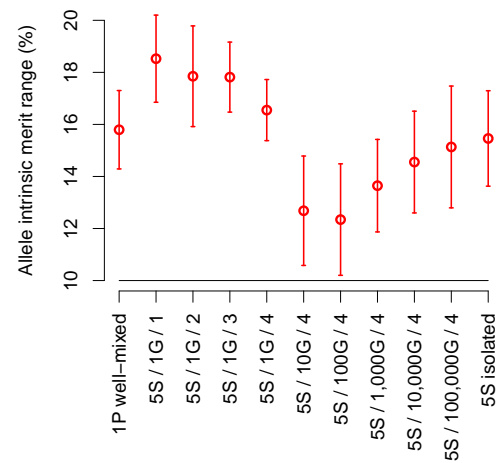


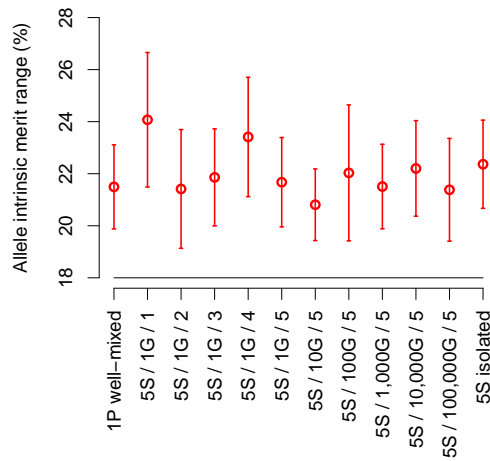
Figure D.1: Number of alleles at the end of the simulation span for different numbers of subpopulations in a star arrangement and different migration rates. The total population size was 1 million.



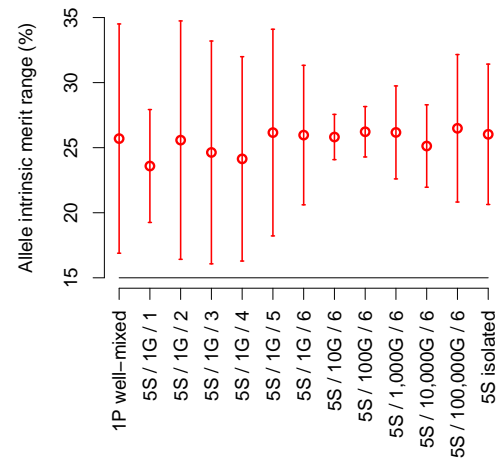
(a) 10000 ind. (total), 5 subpopulations



(b) 100000 ind. (total), 5 subpopulations

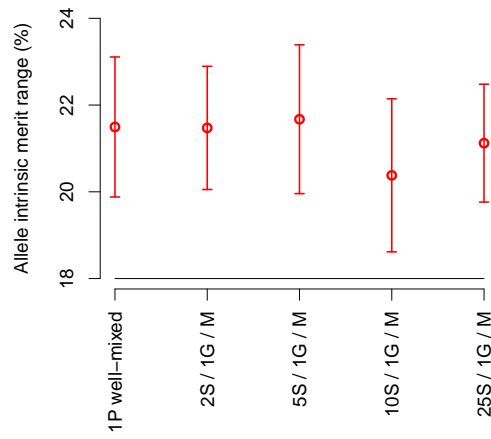


(c) 1 million ind. (total), 5 subpopulations

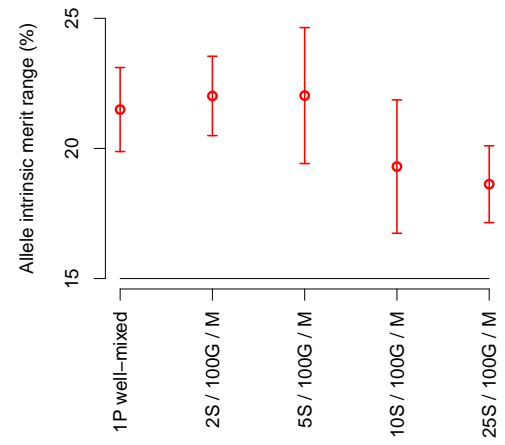


(d) 10 million ind. (total), 5 subpopulations

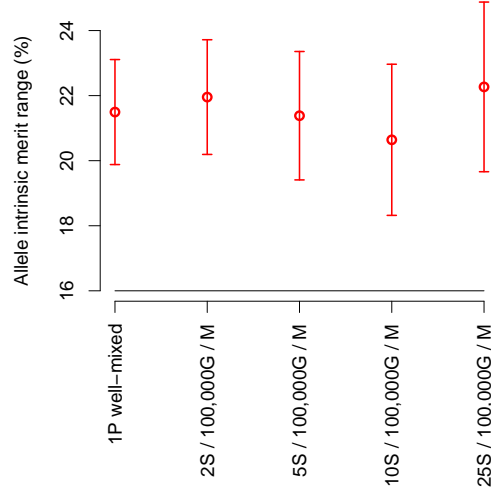
Figure D.2: Range of intrinsic merits of alleles at the end of the simulation span for a metapopulation of 5 subpopulations in a star arrangement and different migration rates and population sizes.



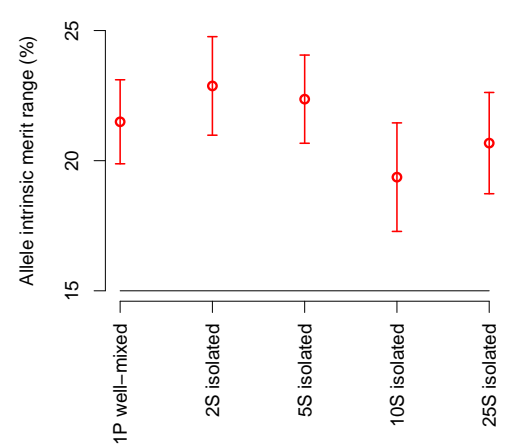
(a) Migration every generation



(b) Migration every 100 generations

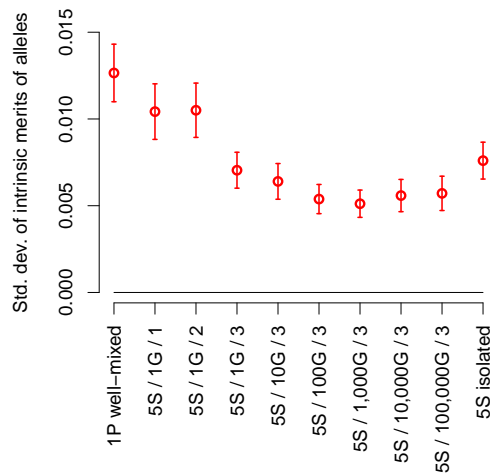


(c) Migration every 100000 generations

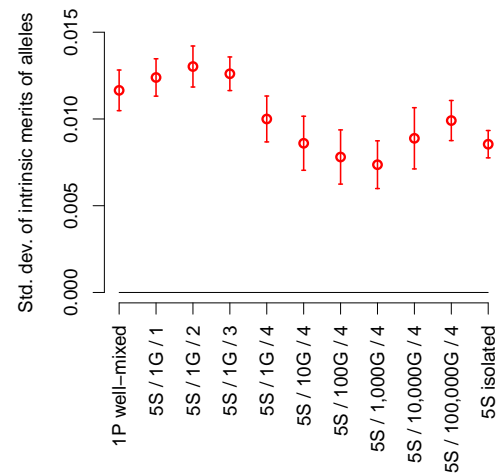


(d) No migration (isolated subpopulations)

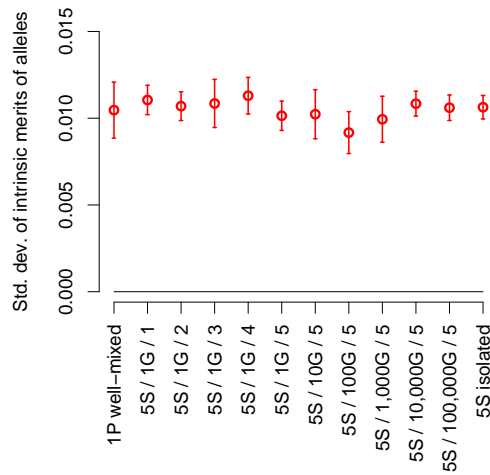
Figure D.3: Range of intrinsic merits of alleles at the end of the simulation span for different numbers of subpopulations in a star arrangement and different migration rates. The total population size was 1 million.



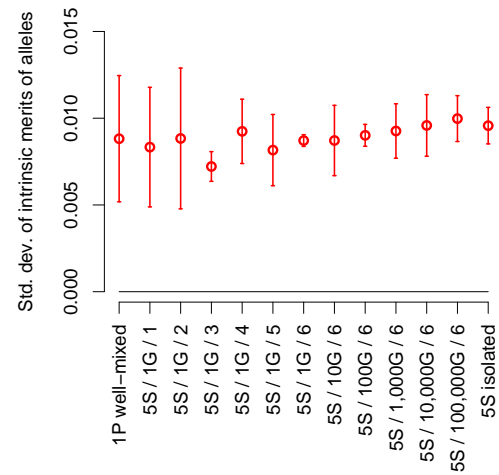
(a) 10000 ind. (total), 5 subpopulations



(b) 100000 ind. (total), 5 subpopulations

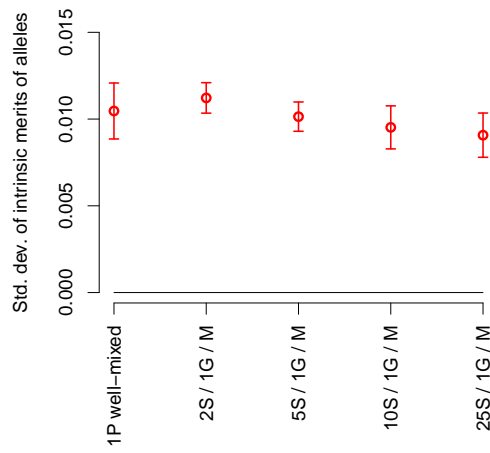


(c) 1 million ind. (total), 5 subpopulations

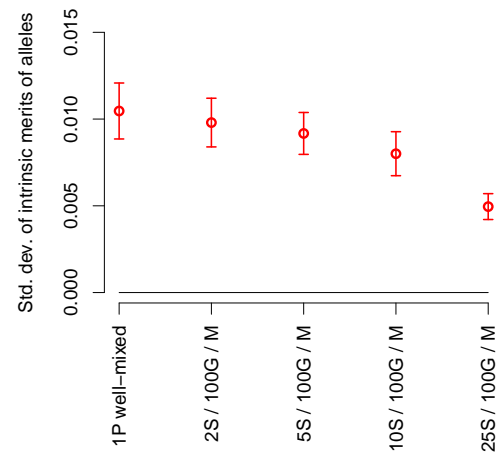


(d) 10 million ind. (total), 5 subpopulations

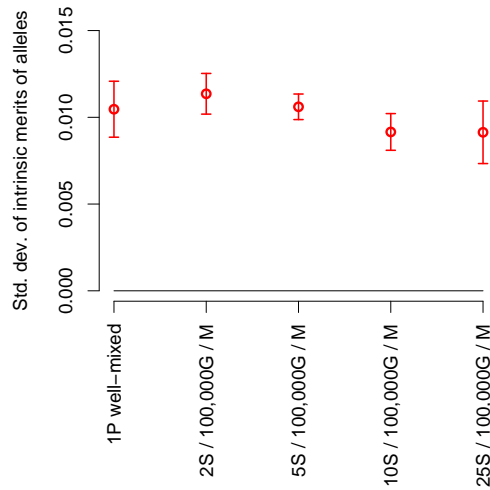
Figure D.4: Weighted standard deviation in allele intrinsic merit at the end of the simulation span for a metapopulation of 5 subpopulations in a star arrangement and different migration rates and population sizes.



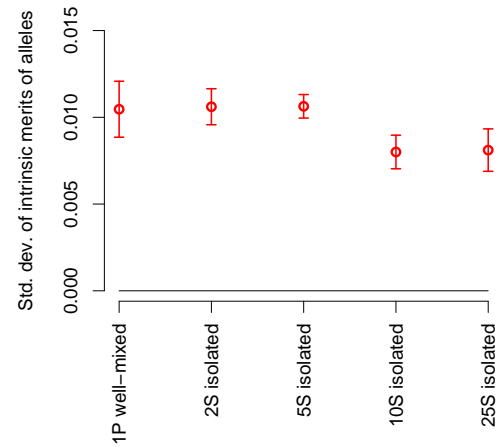
(a) Migration every generation



(b) Migration every 100 generations



(c) Migration every 100000 generations



(d) No migration (isolated subpopulations)

Figure D.5: Weighted standard deviation in allele intrinsic merit at the end of the simulation span for different numbers of subpopulations in a star arrangement and different migration rates. The total population size was 1 million.

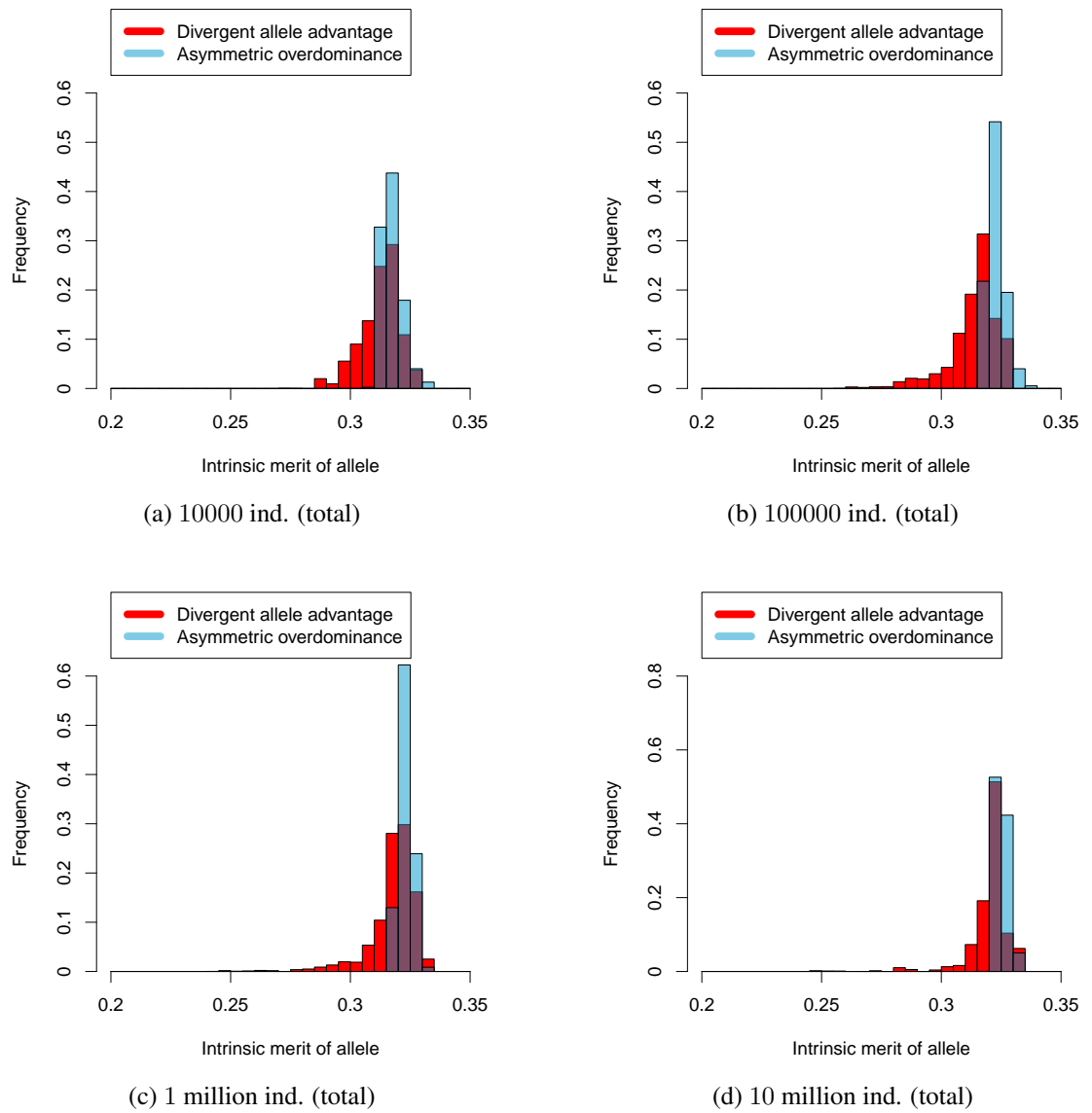


Figure D.6: Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every generation.

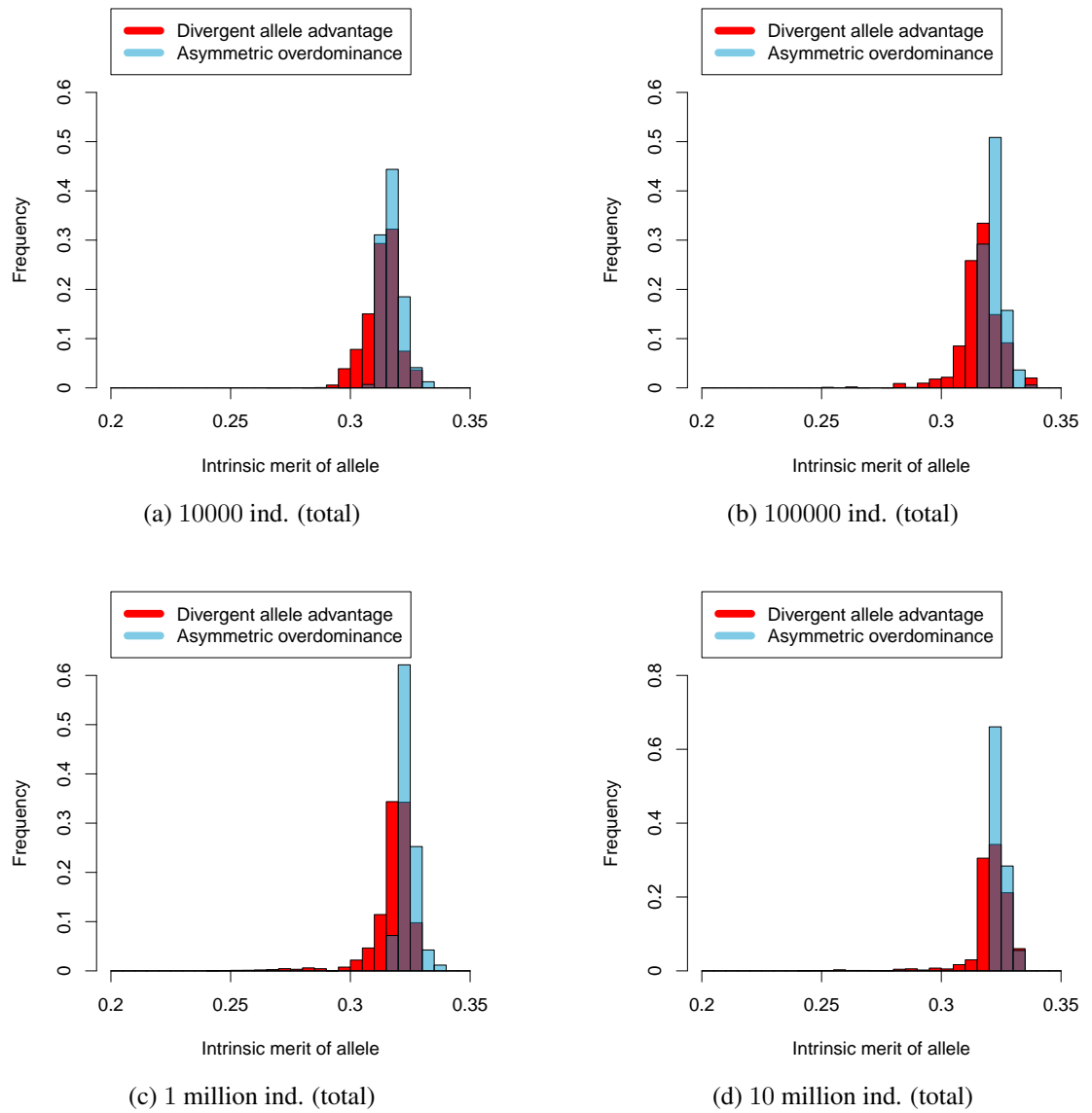


Figure D.7: Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every 100 generations.

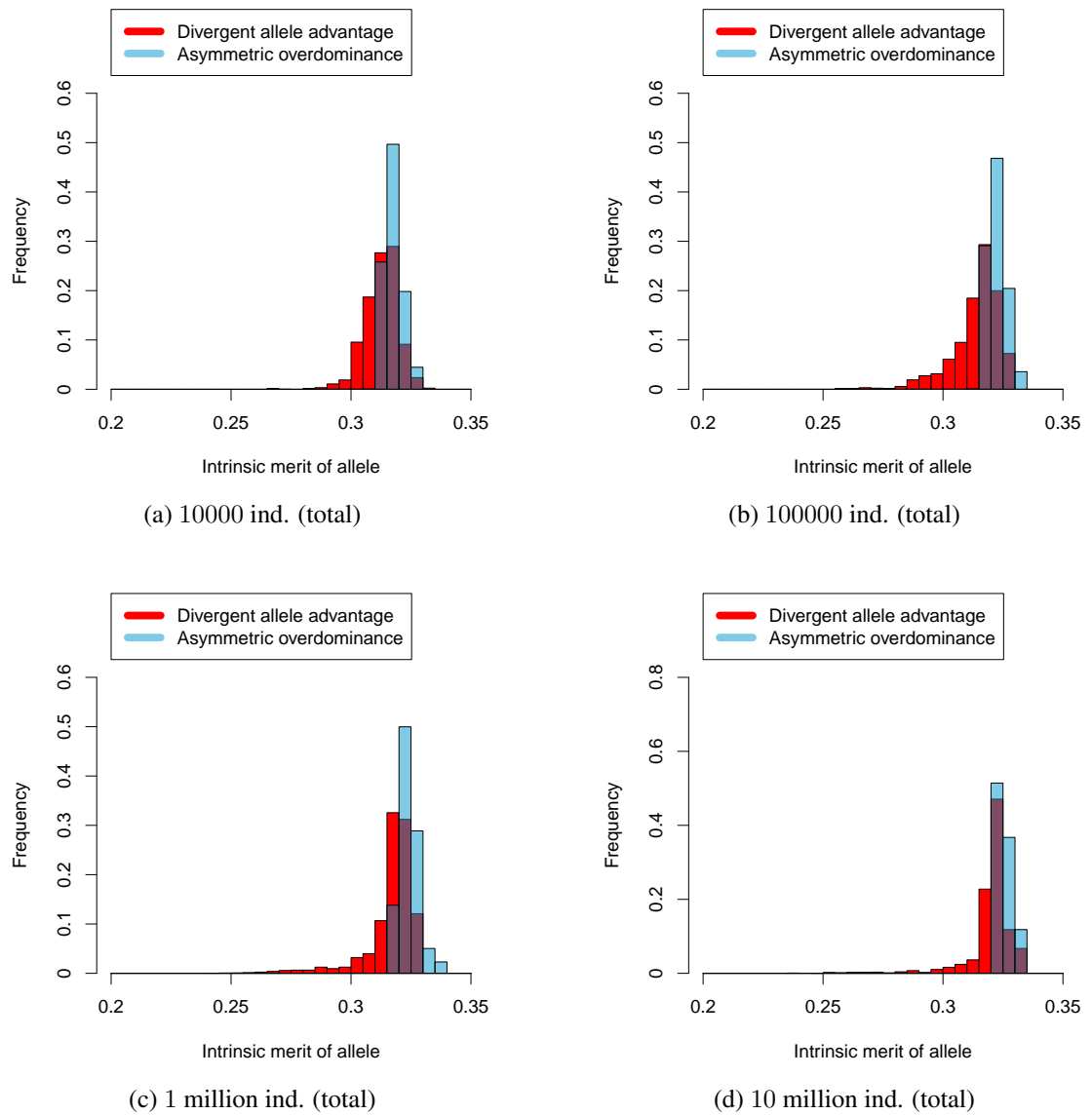


Figure D.8: Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every 100000 generations.

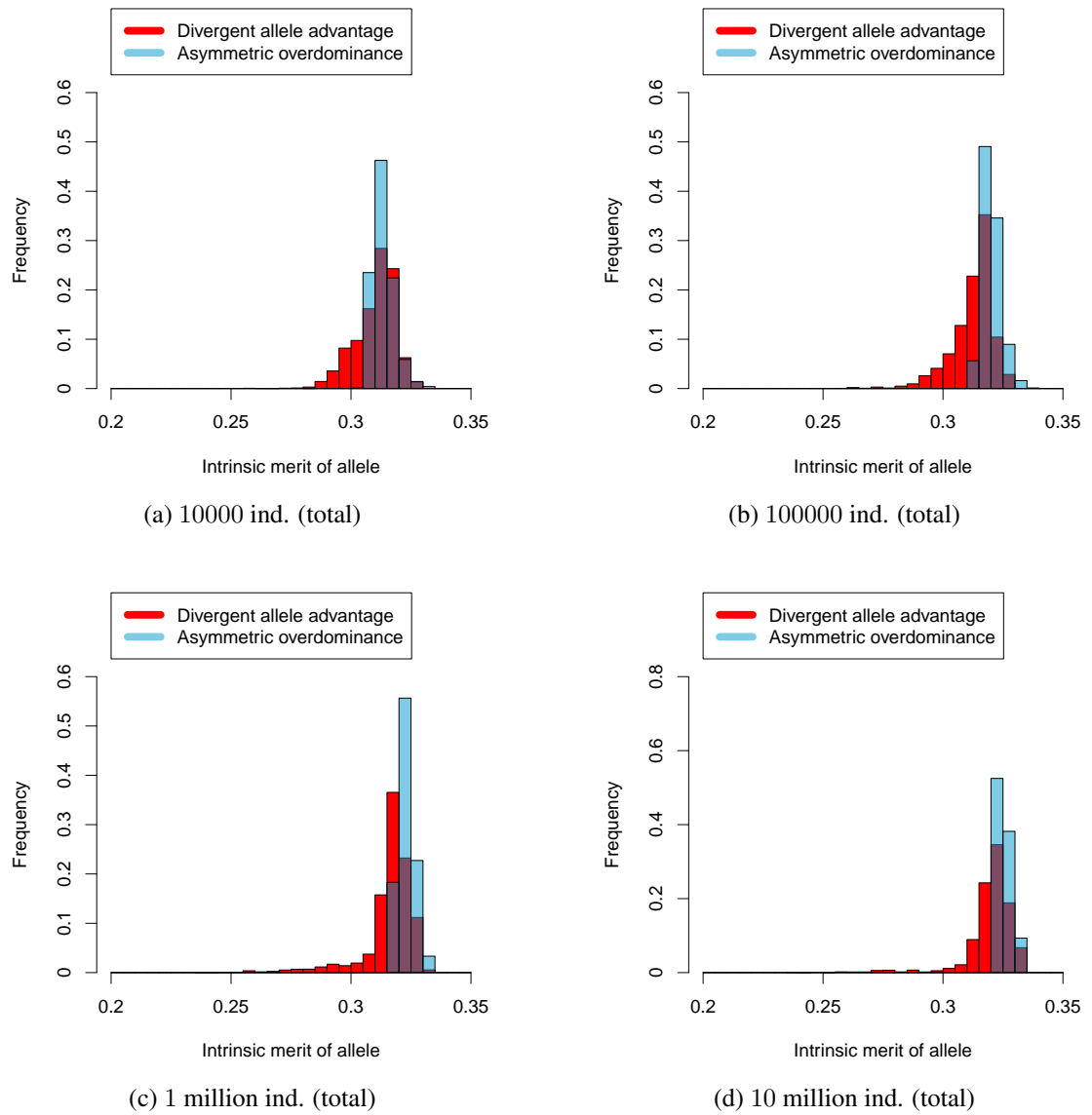


Figure D.9: Distribution of intrinsic merit of alleles for a metapopulation in a star arrangement, all population sizes tested and a metapopulation of 5 isolated subpopulations.

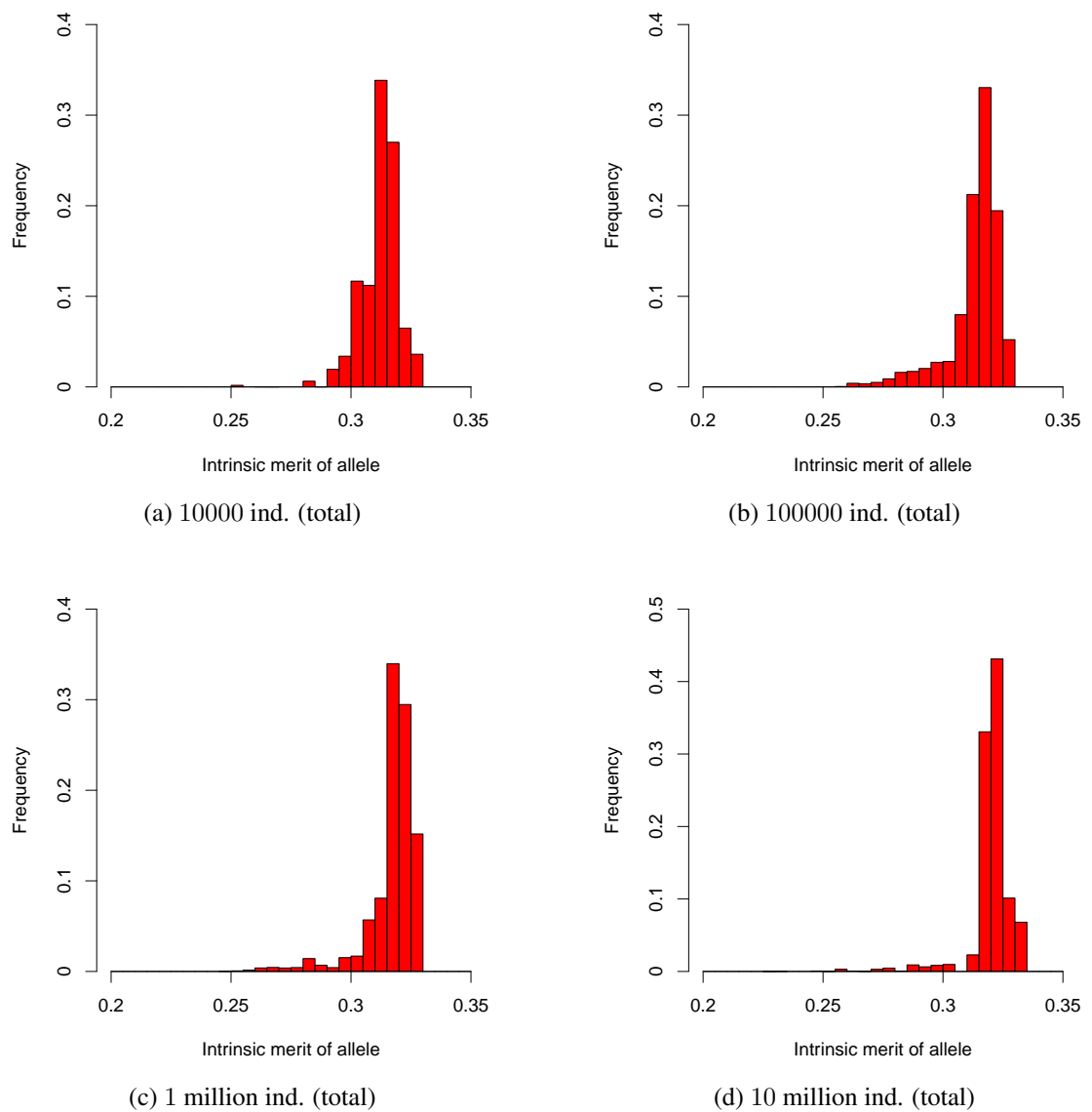


Figure D.10: Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every generation.

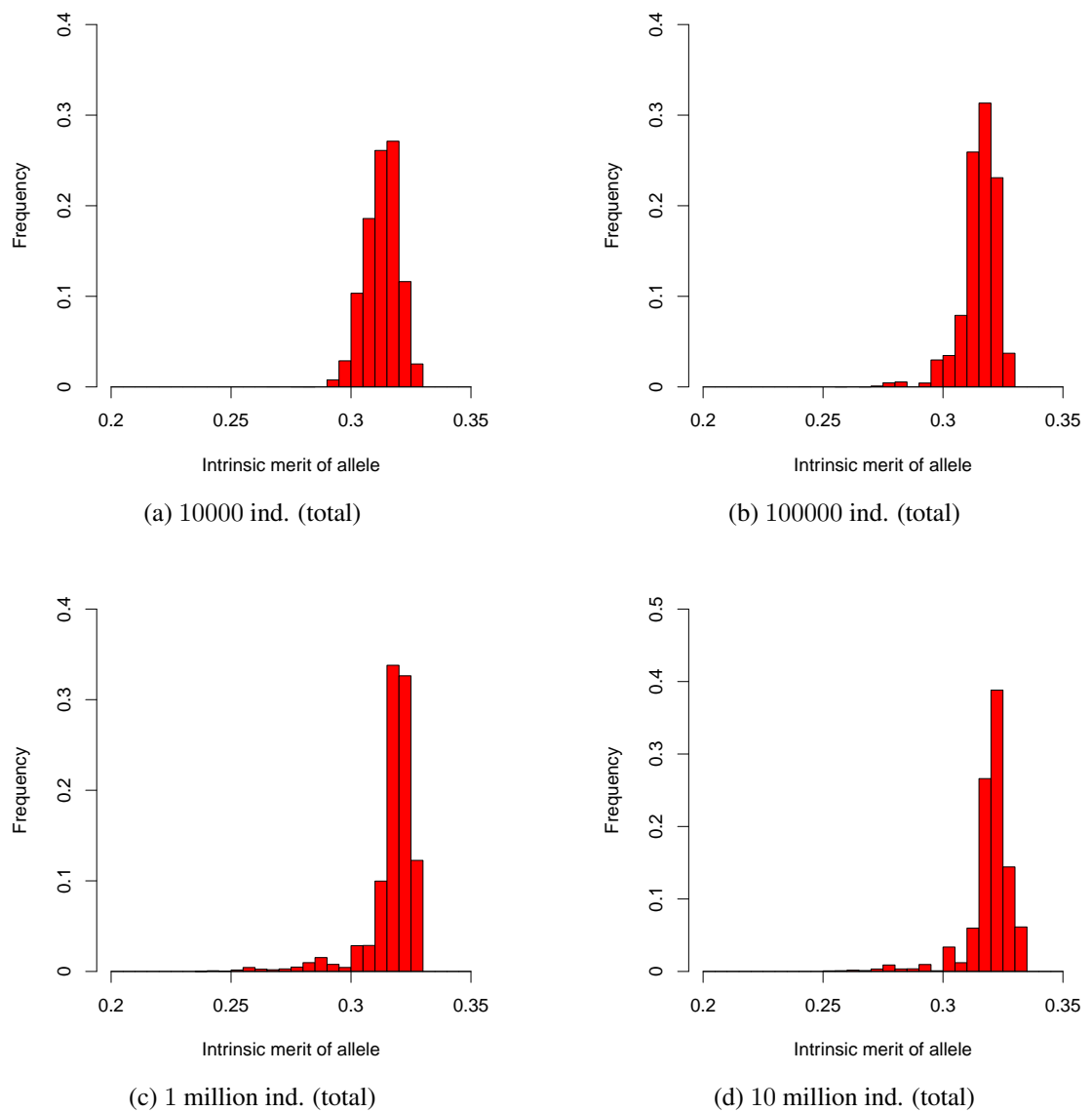


Figure D.11: Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every 100 generations.

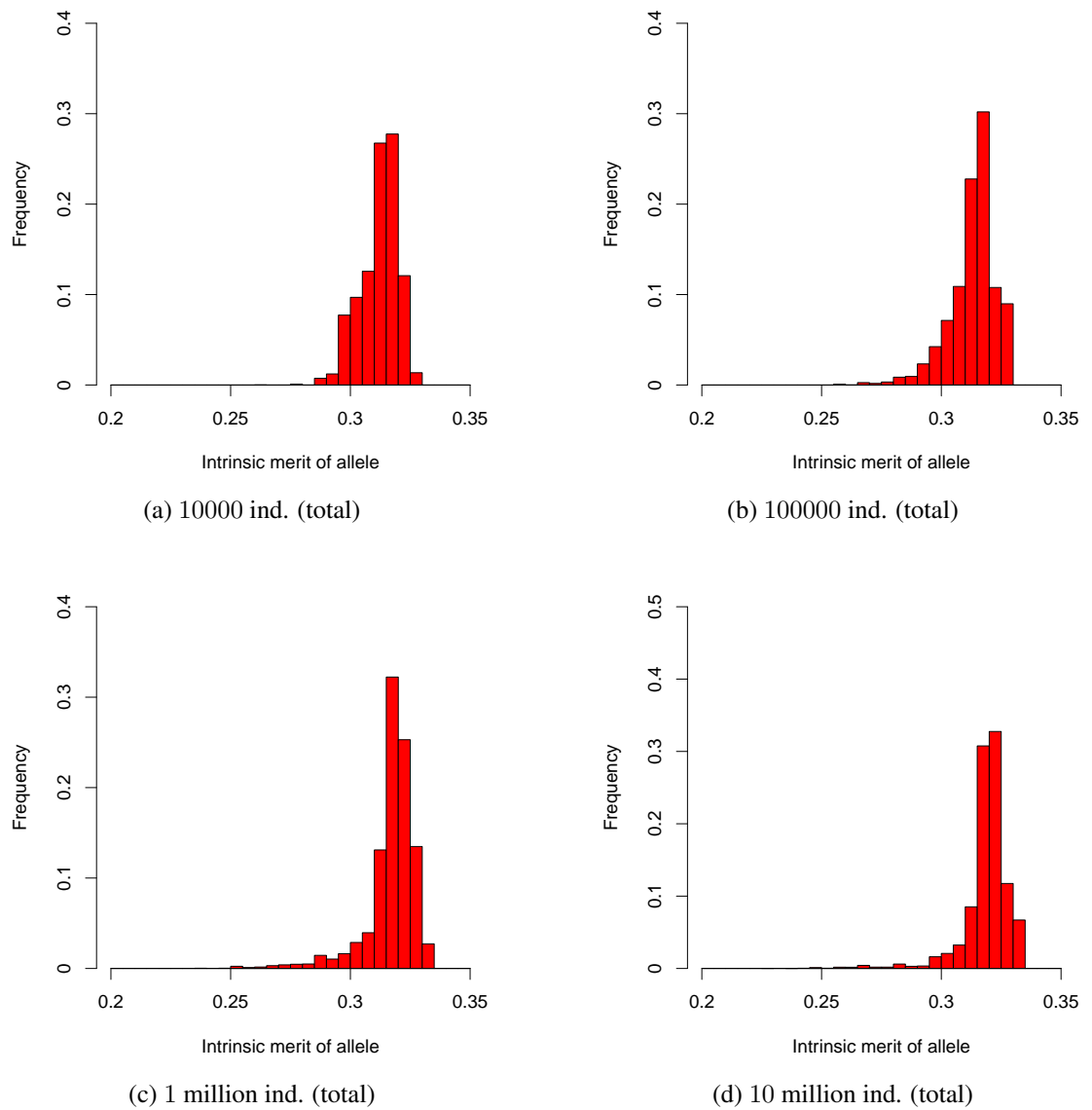
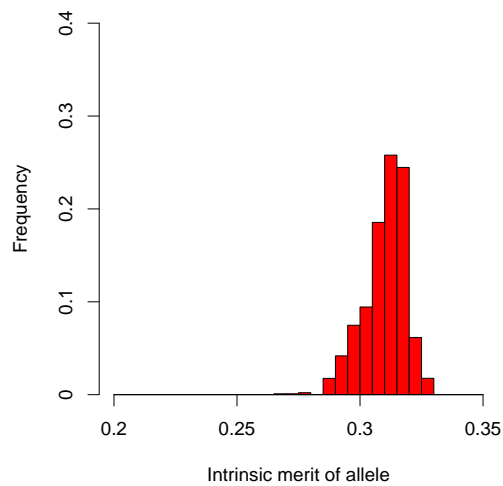
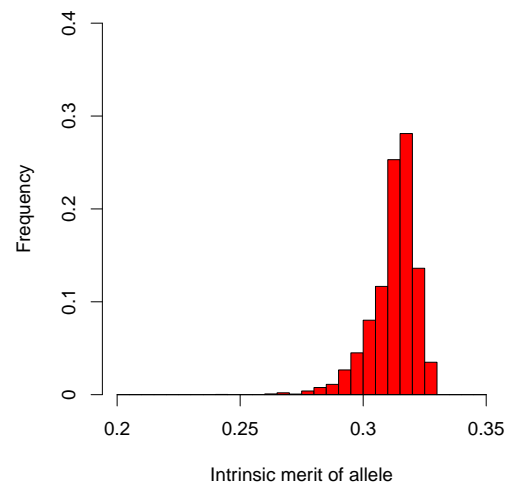


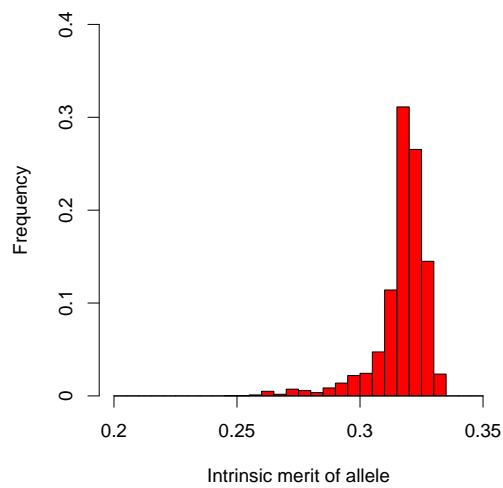
Figure D.12: Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every 100000 generations.



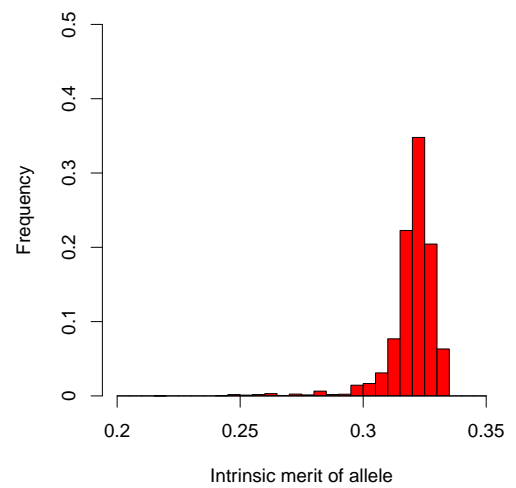
(a) 10000 ind. (total)



(b) 100000 ind. (total)

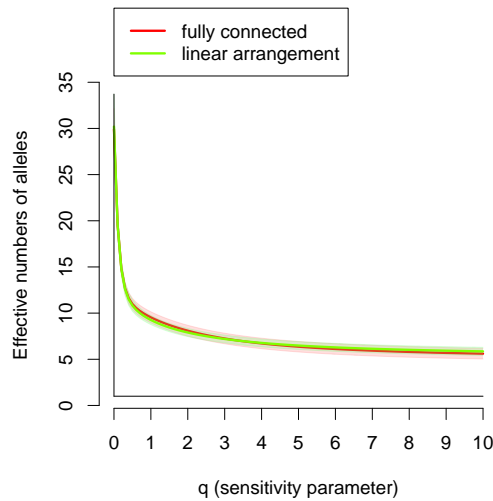


(c) 1 million ind. (total)

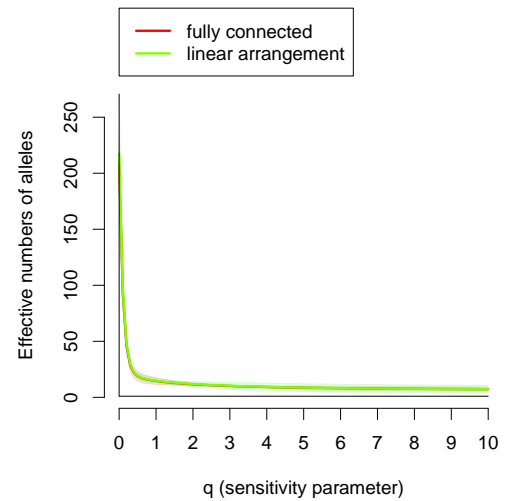


(d) 10 million ind. (total)

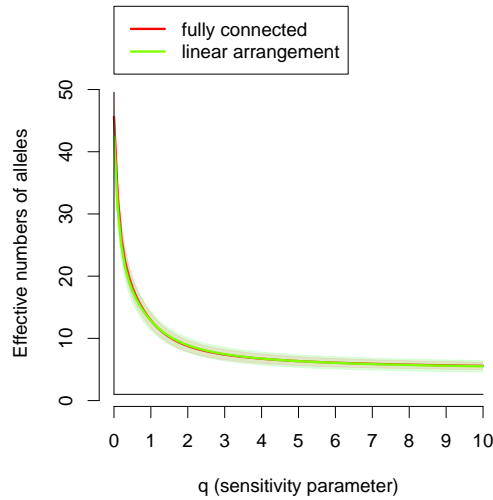
Figure D.13: Distribution of intrinsic merit of alleles for a metapopulation in a linear arrangement, all population sizes tested and a metapopulation of 5 isolated subpopulations.



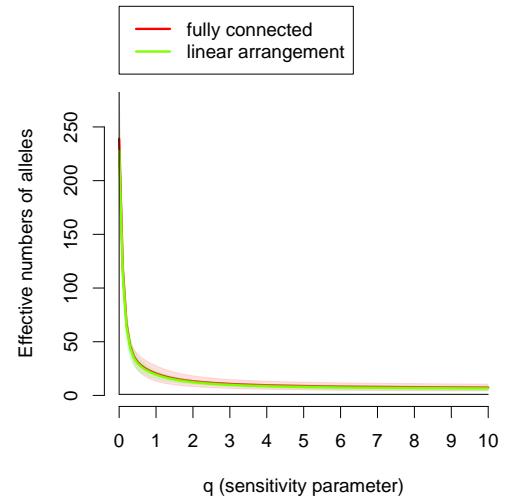
(a) 1 million ind. (total), migration every generation



(b) 10 million ind. (total), migration every generation



(c) 1 million ind. (total), migration every 100000 generations



(d) 10 million ind. (total), migration every 100000 generations

Figure D.14: Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations in a star arrangement and a metapopulation of 5 subpopulations in a linear arrangement, for a population size of 1 million (left) and 10 million (right) respectively.

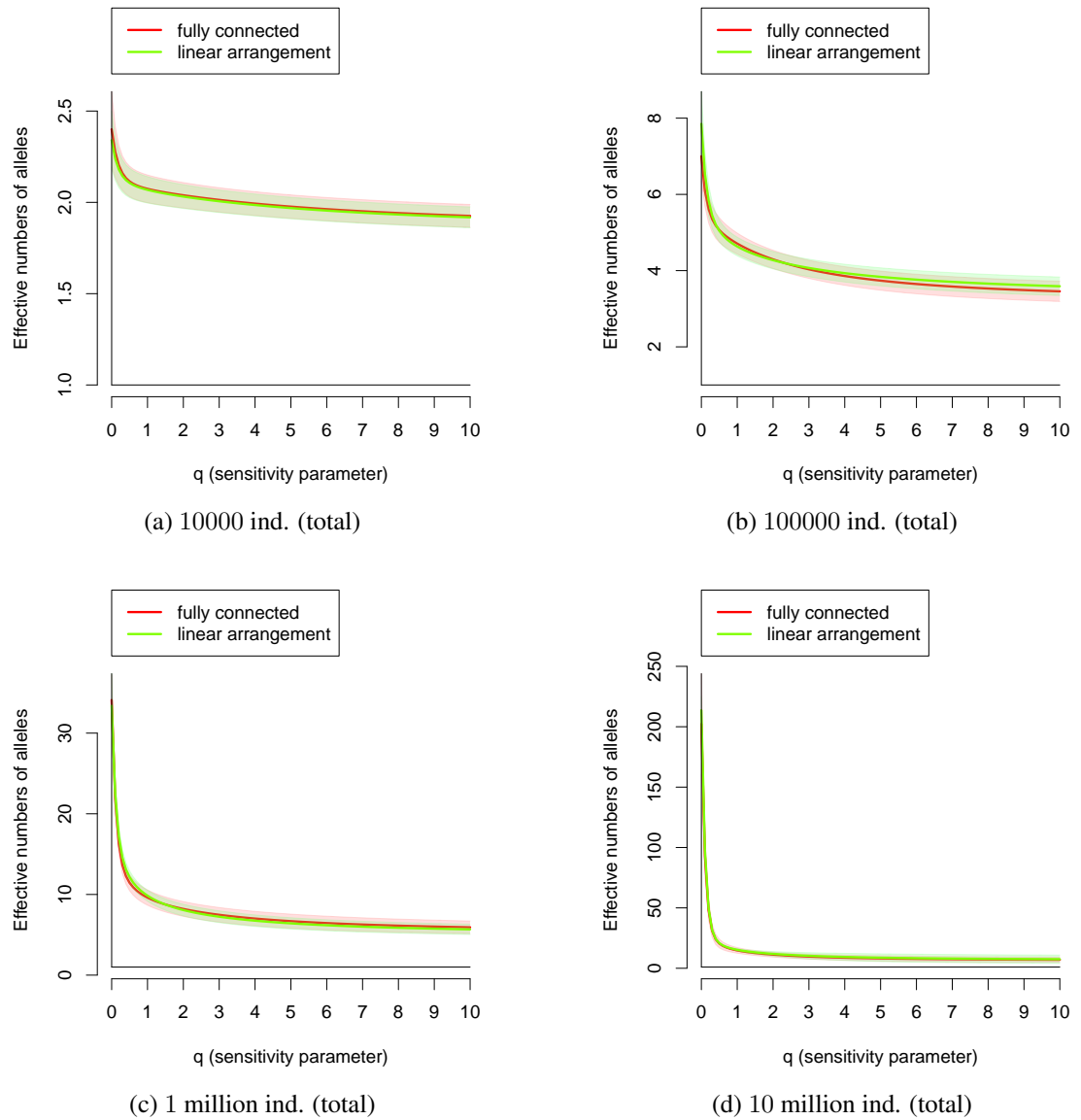


Figure D.15: Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations in a star arrangement and a metapopulation of 5 subpopulations in a linear arrangement, for all population sizes tested and a migration rate of 2 individuals every 100 generations.

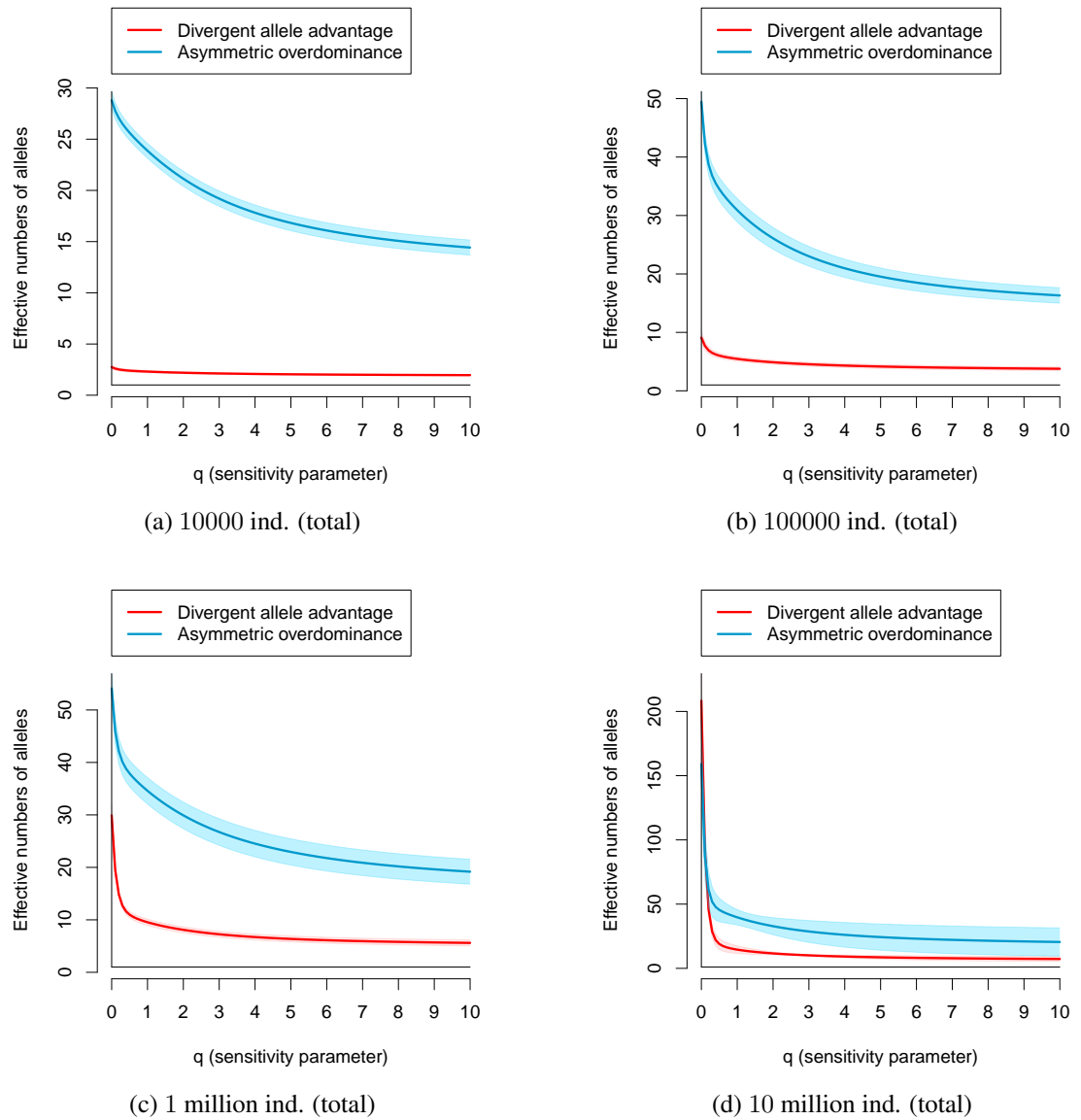


Figure D.16: Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations in a star arrangement, the DAA and asymmetric overdominance models and a migration rate of 2 individuals every generation.

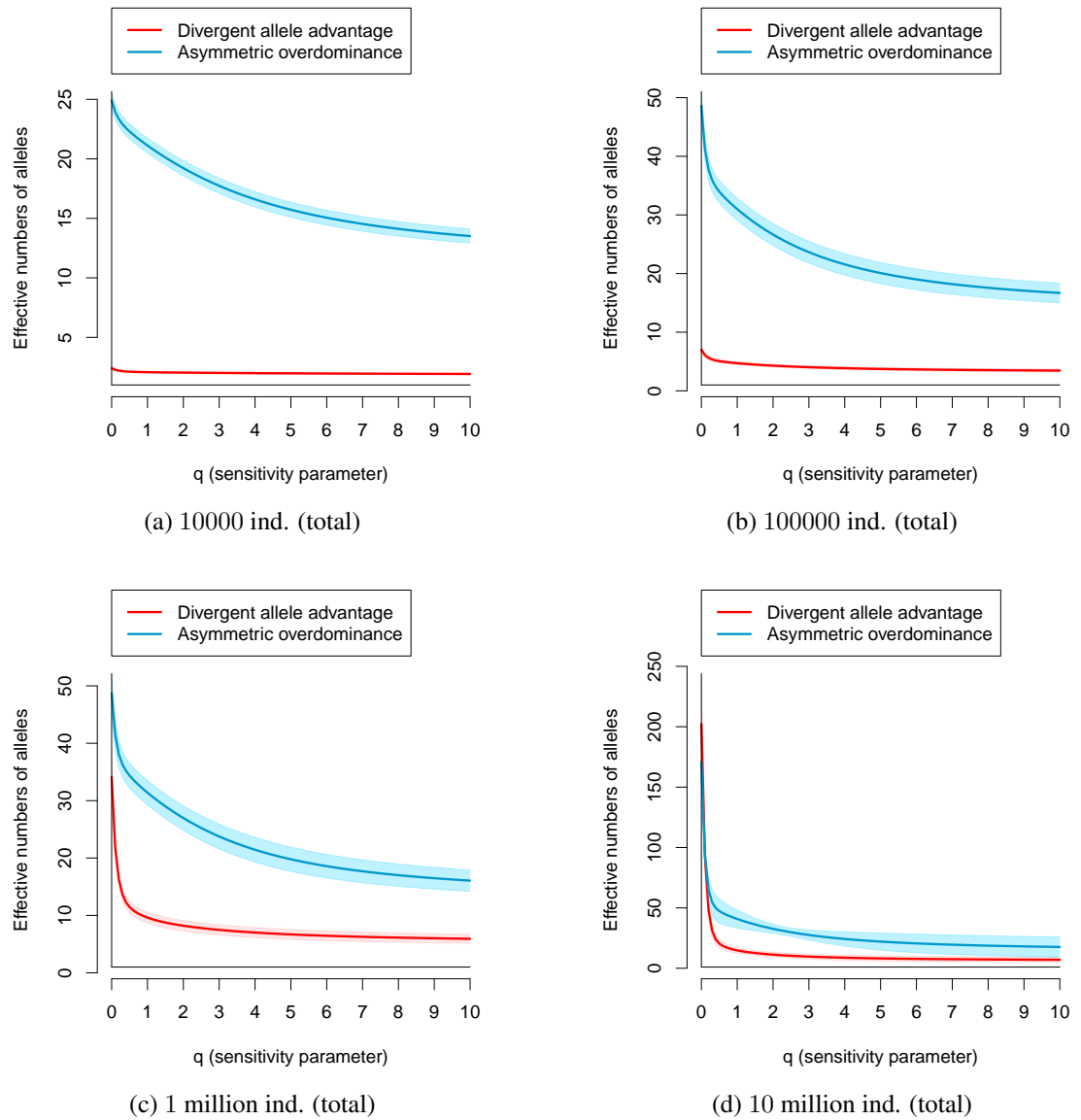


Figure D.17: Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations in a star arrangement, the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 100 generations.

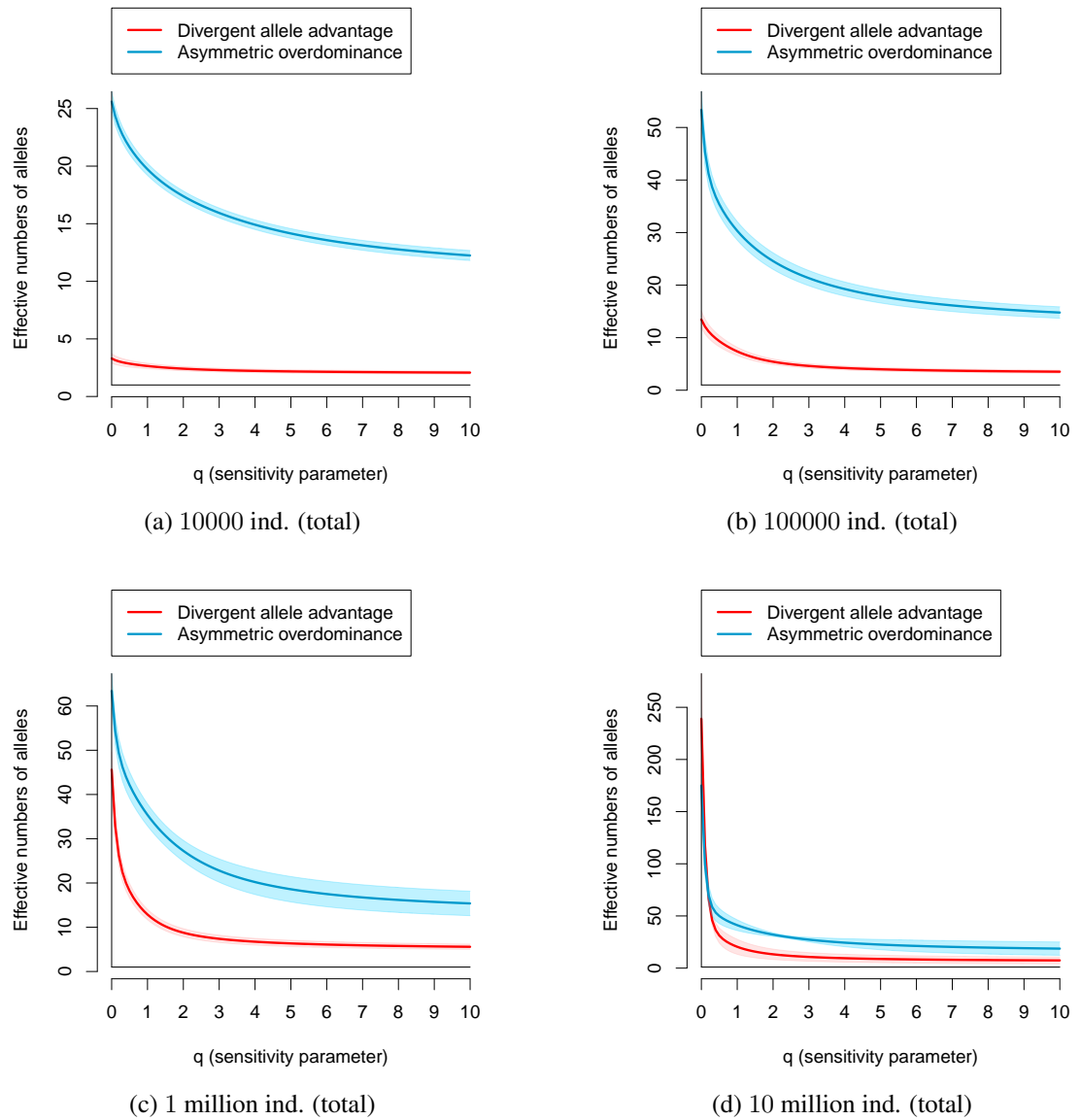


Figure D.18: Diversity profiles (naive diversity) for a metapopulation of 5 subpopulations in a star arrangement, the DAA and asymmetric overdominance models and a migration rate of 2 individuals every 100000 generations.

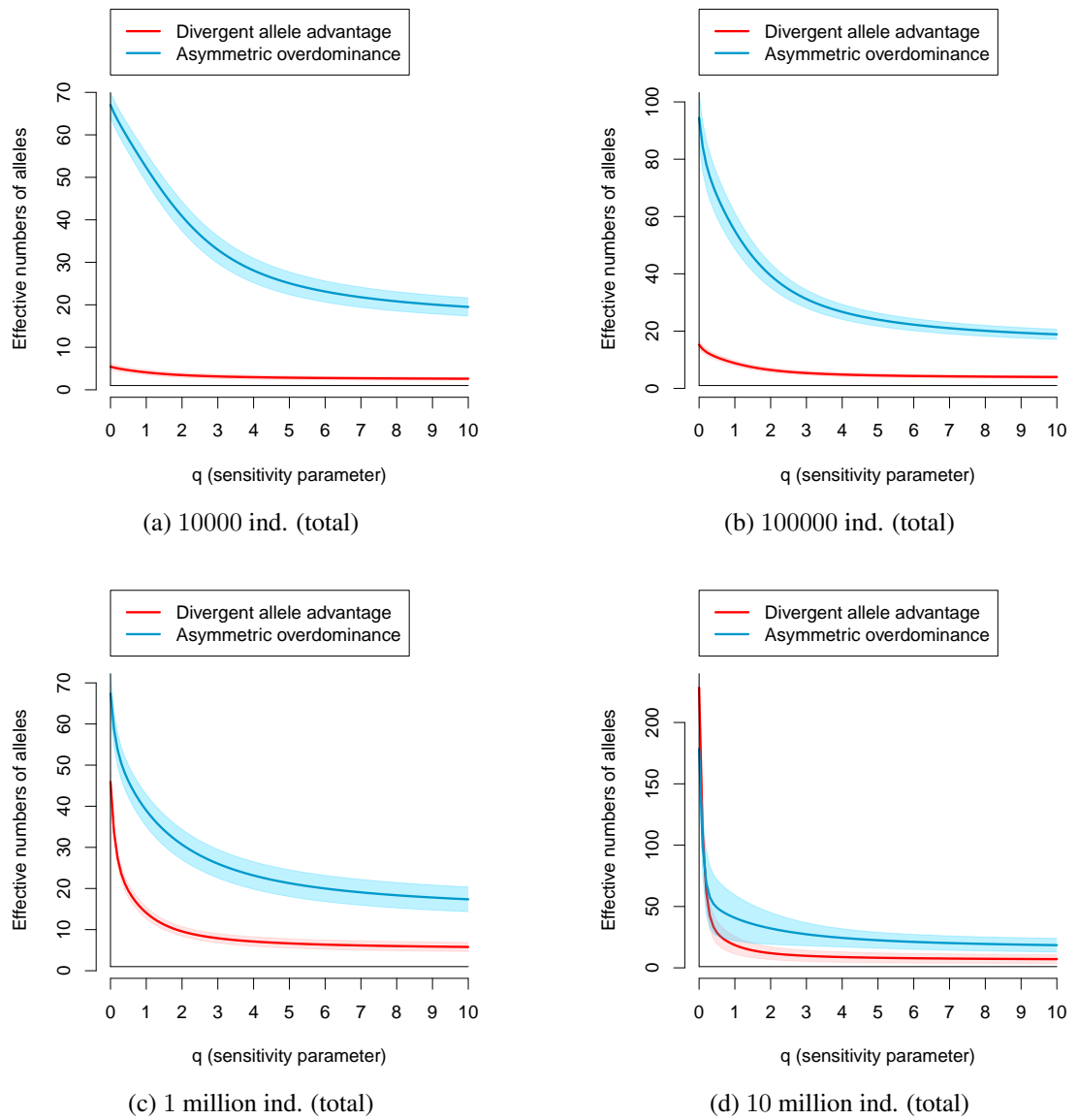


Figure D.19: Diversity profiles (naive diversity) for a metapopulation of 5 isolated subpopulations and the DAA and asymmetric overdominance models.

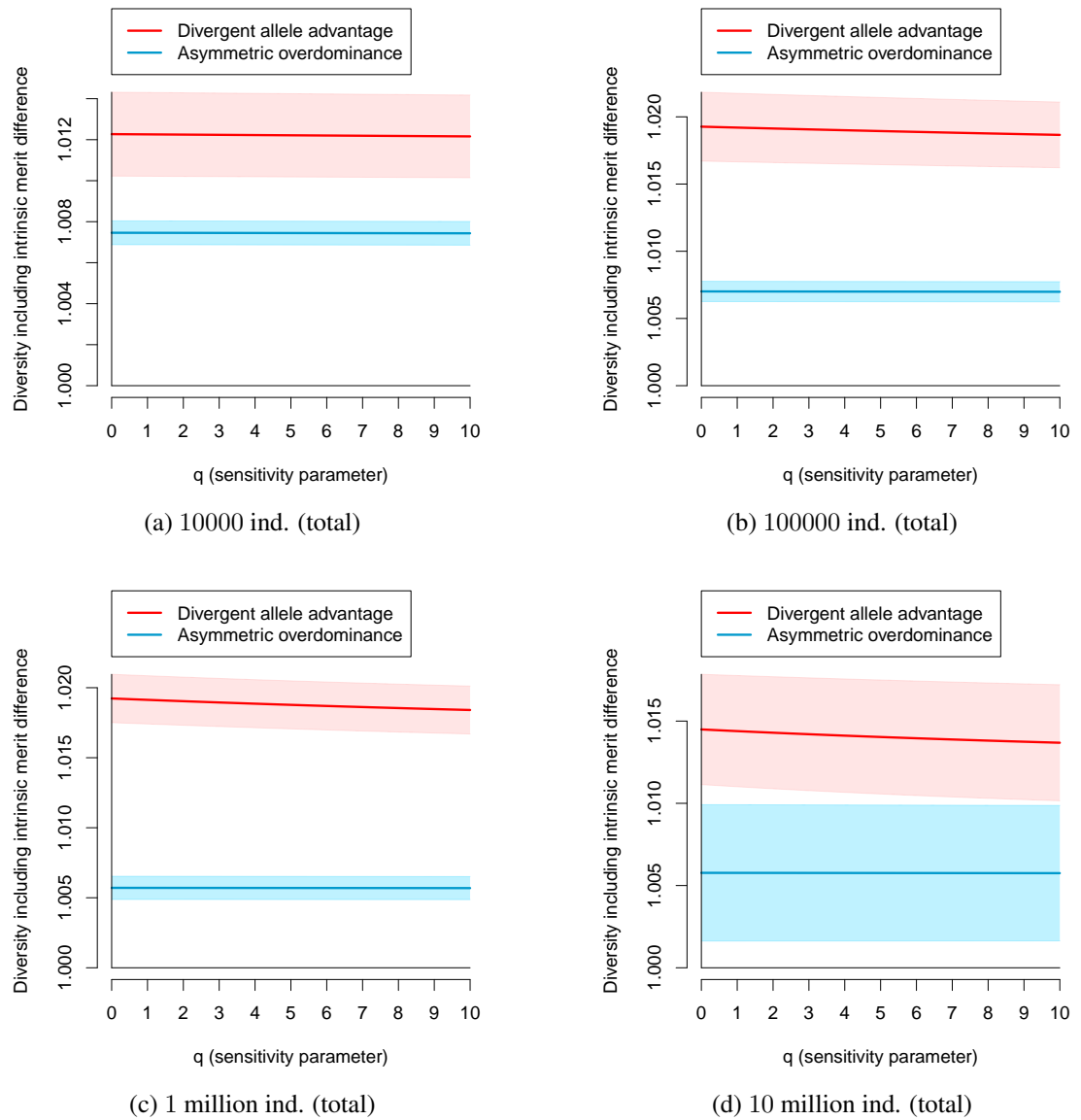


Figure D.20: Diversity profiles (similarity-sensitive diversity based on intrinsic merit difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every generation.

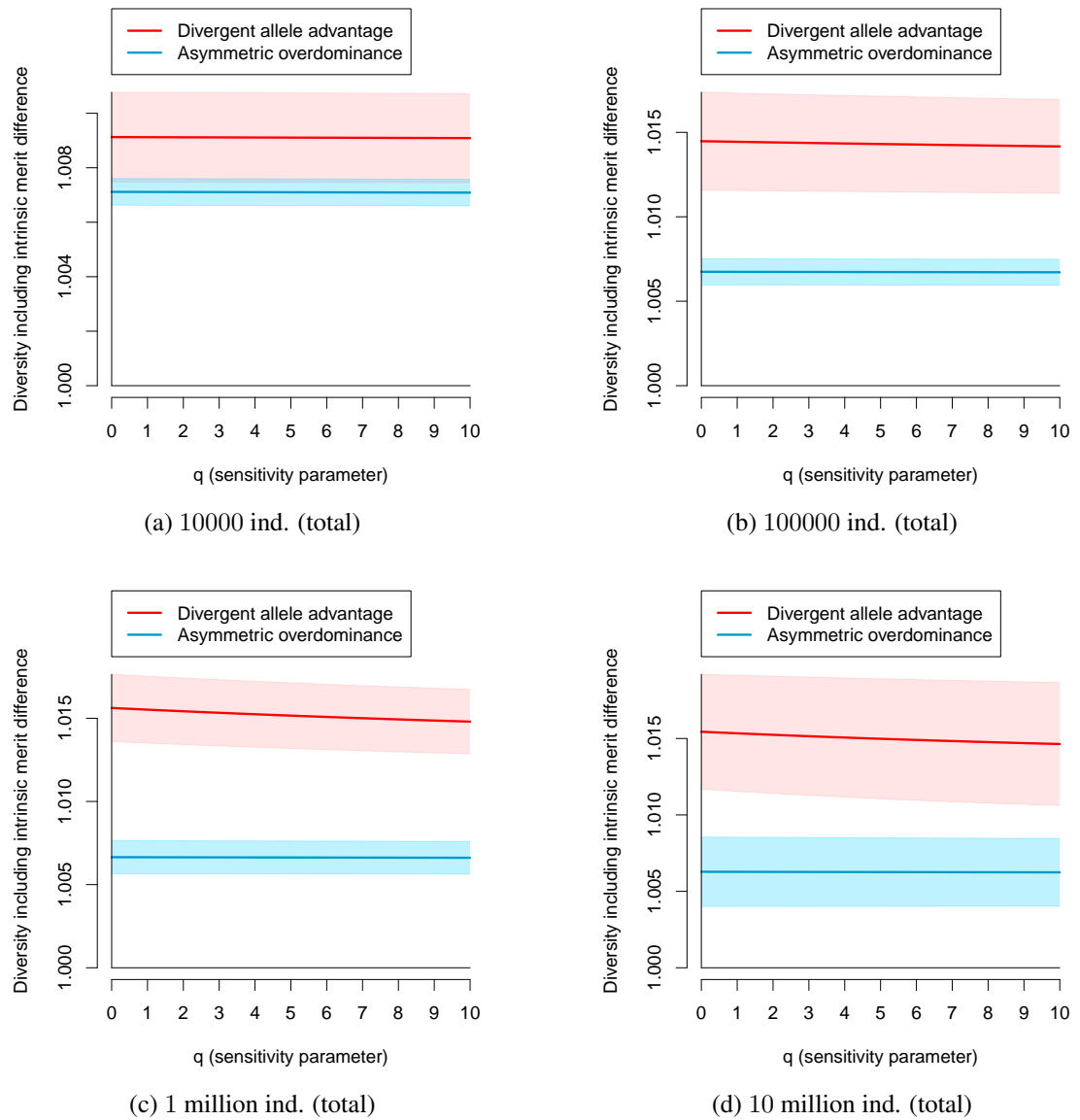


Figure D.21: Diversity profiles (similarity-sensitive diversity based on intrinsic merit difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every 100 generations.

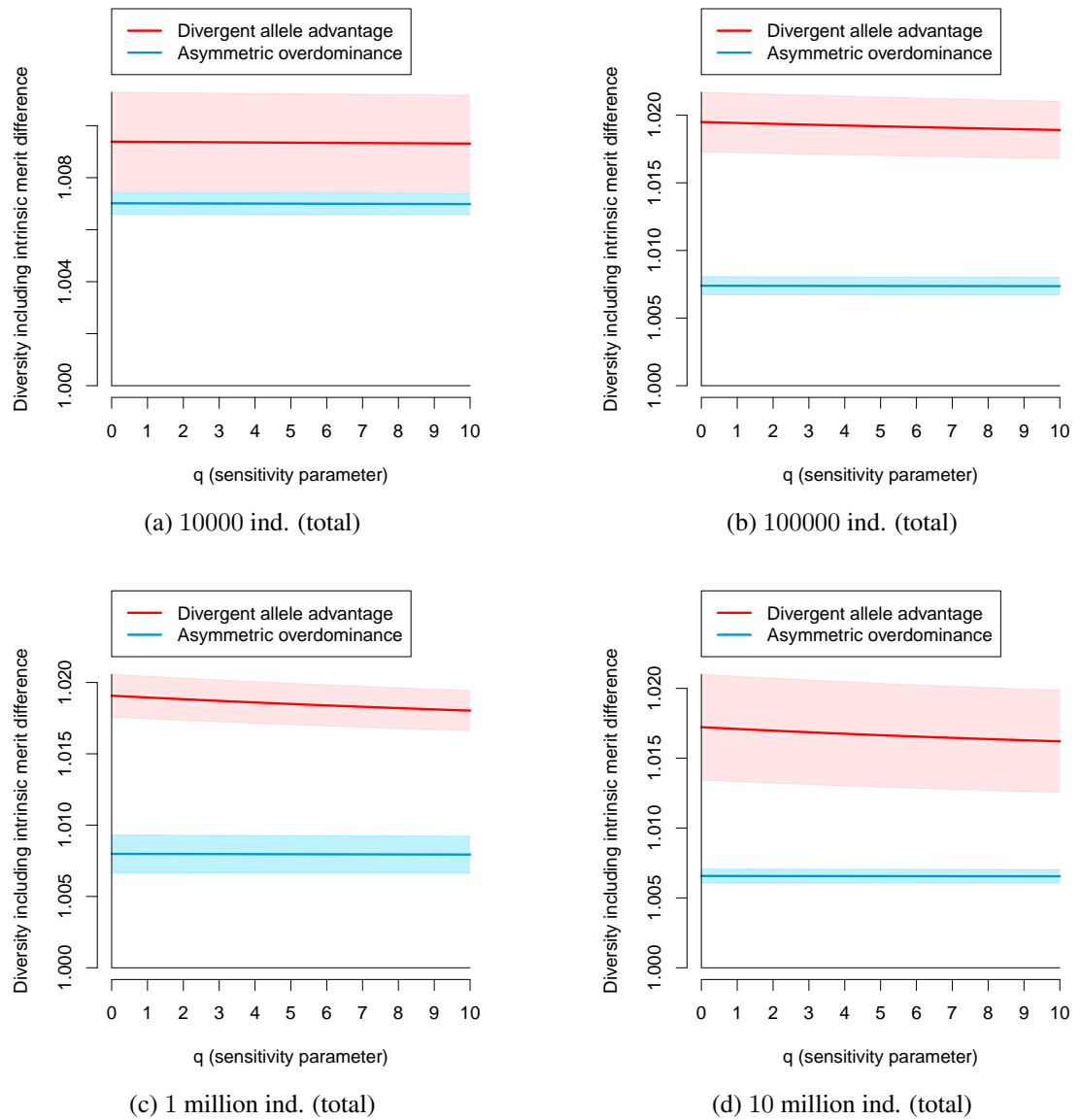


Figure D.22: Diversity profiles (similarity-sensitive diversity based on intrinsic merit difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every 100000 generations.

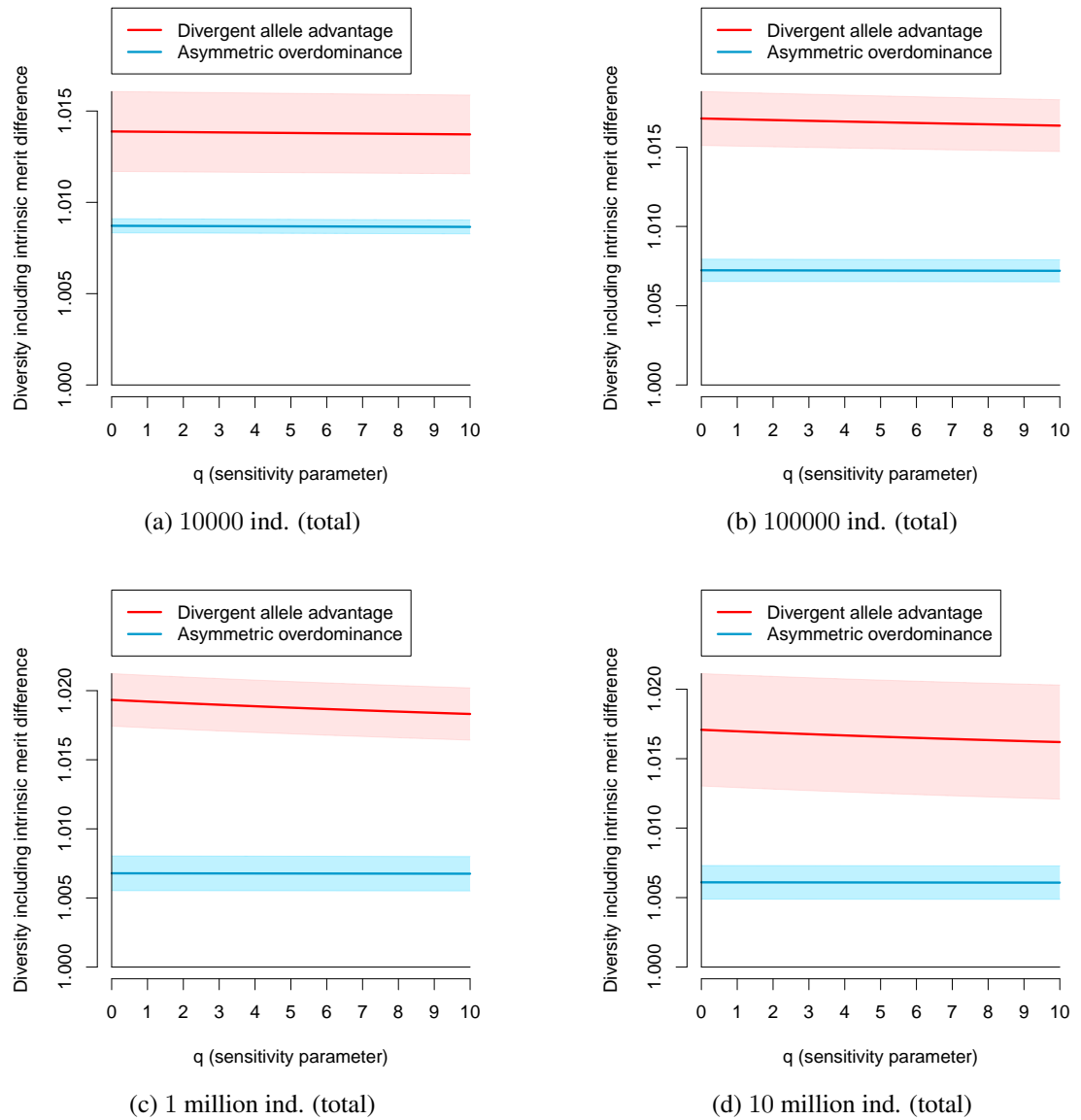


Figure D.23: Diversity profiles (similarity-sensitive diversity based on intrinsic merit difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 isolated subpopulations.

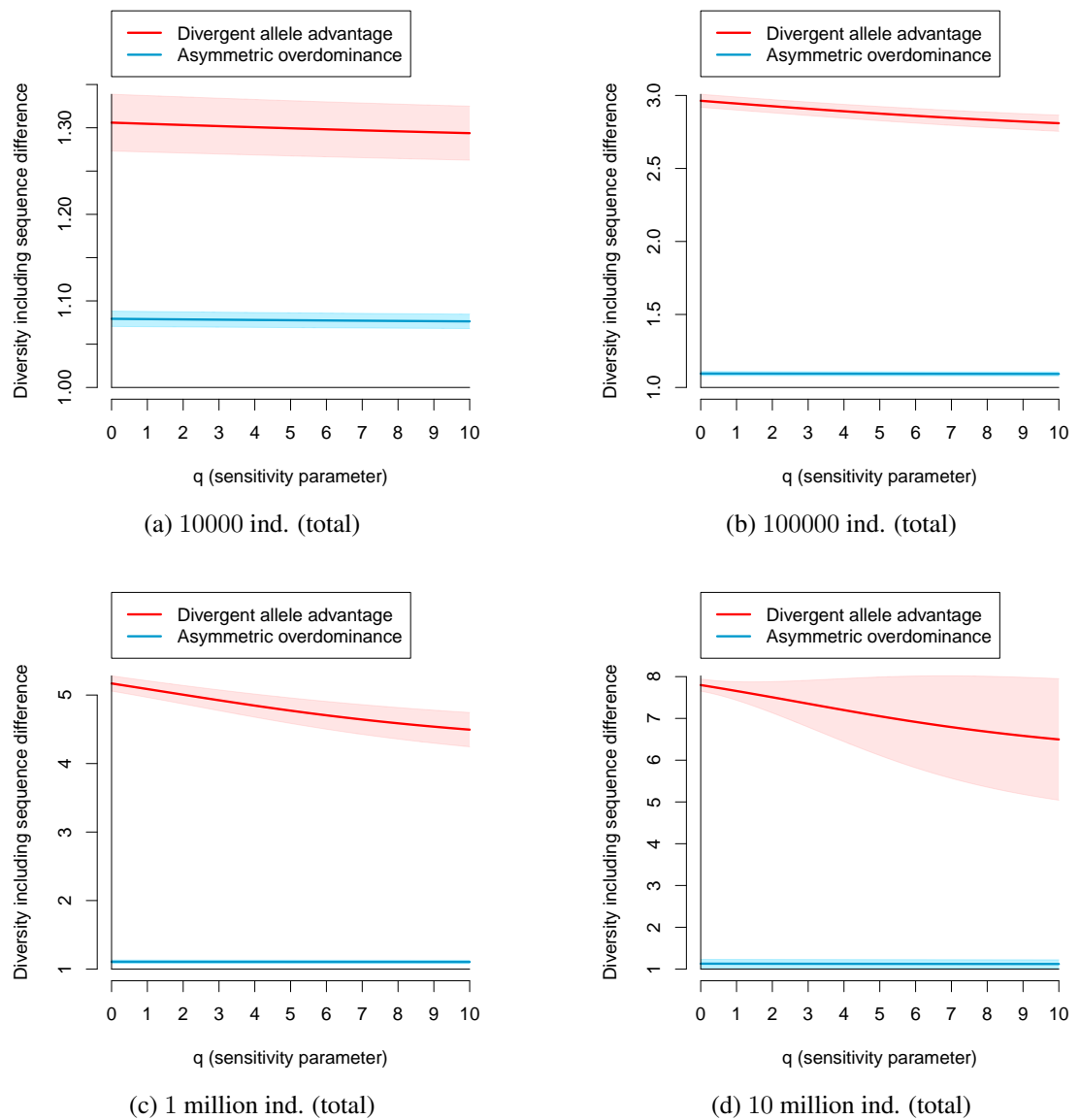


Figure D.24: Diversity profiles (similarity-sensitive diversity based on sequence difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every generation.

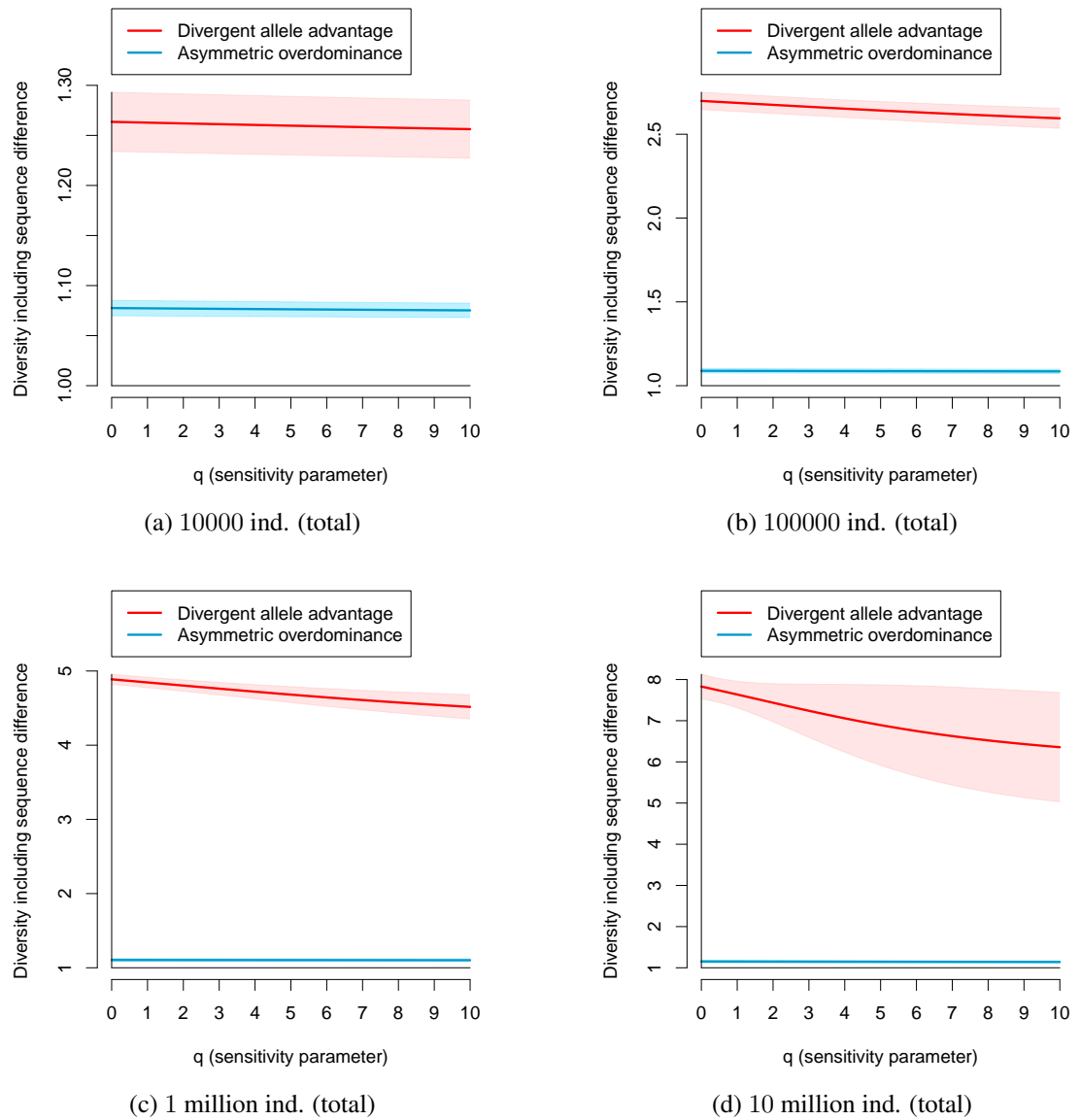


Figure D.25: Diversity profiles (similarity-sensitive diversity based on sequence difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every 100 generations.

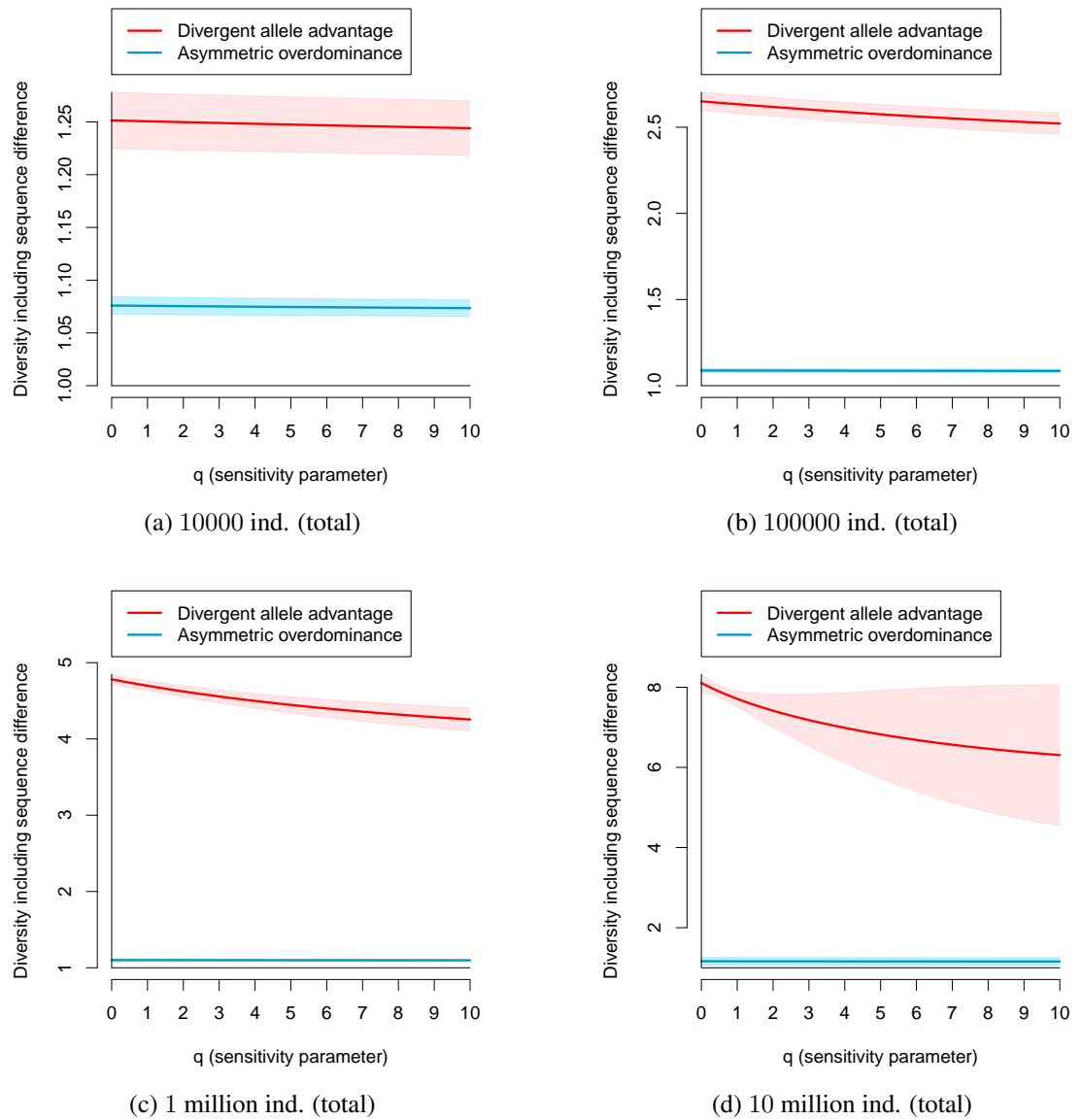


Figure D.26: Diversity profiles (similarity-sensitive diversity based on sequence difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 subpopulations in a star arrangement and a migration rate of 2 individuals every 100,000 generations.

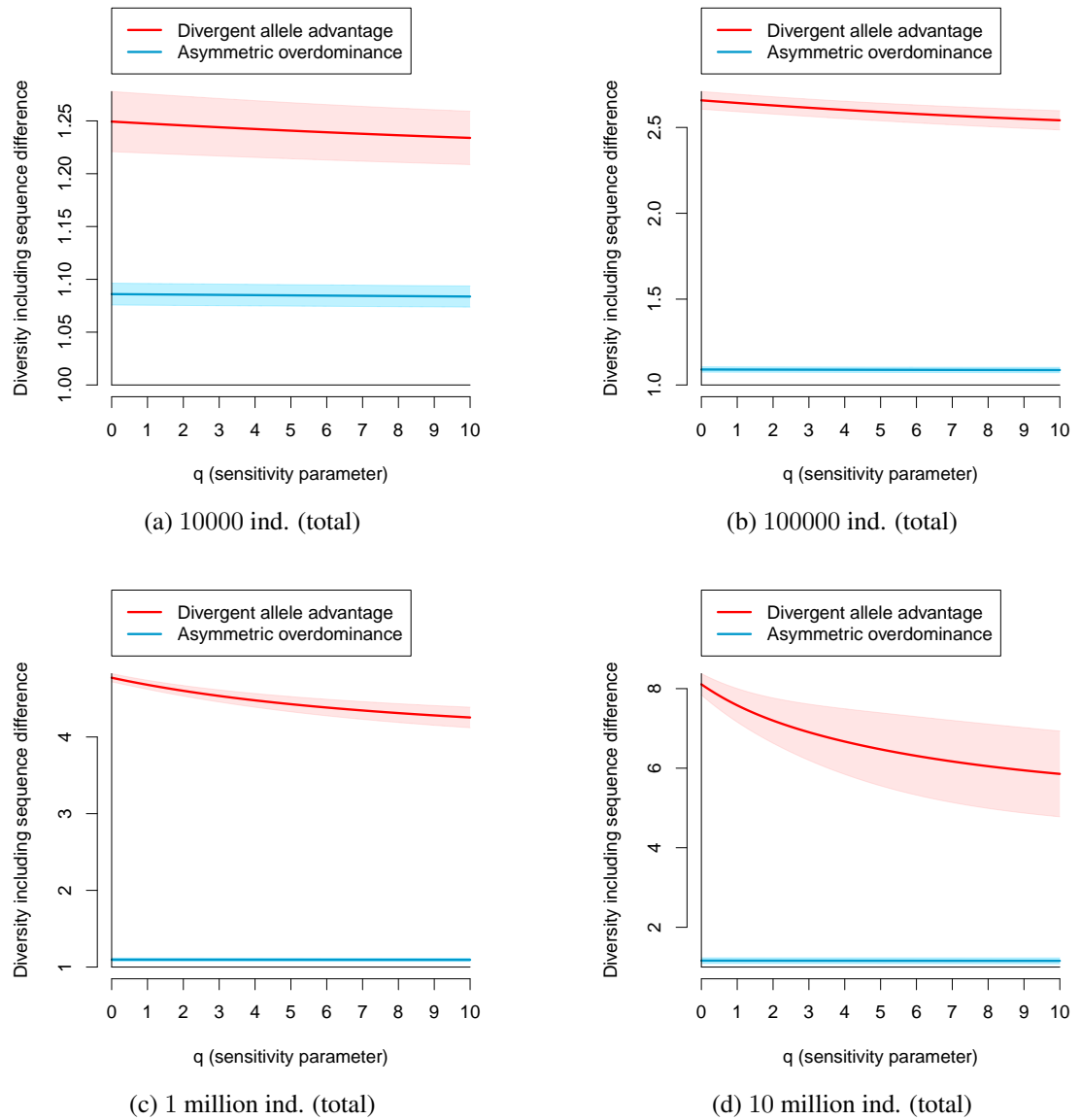
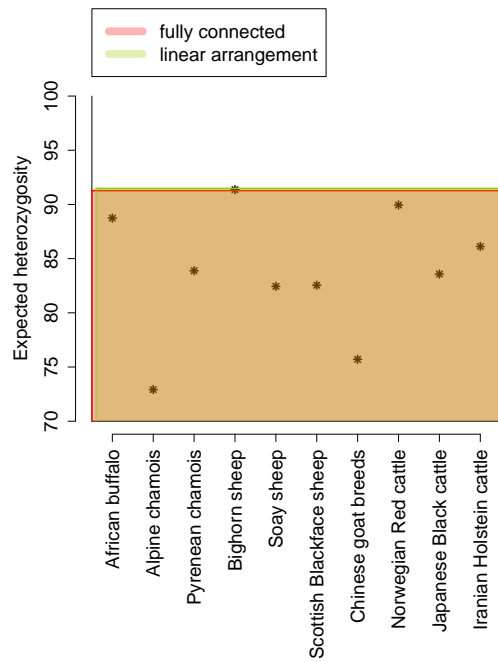
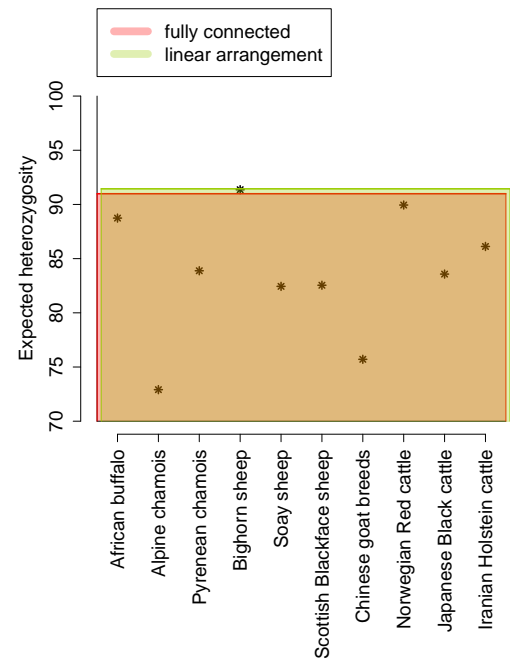


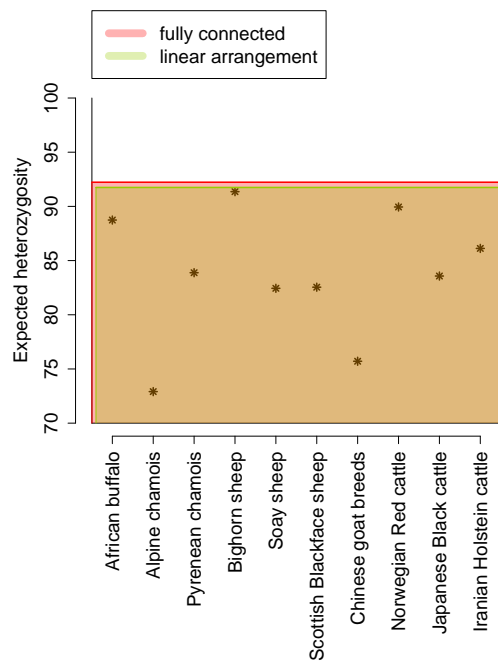
Figure D.27: Diversity profiles (similarity-sensitive diversity based on sequence difference between alleles) – comparison between the DAA and asymmetric overdominance models for a metapopulation of 5 isolated subpopulations.



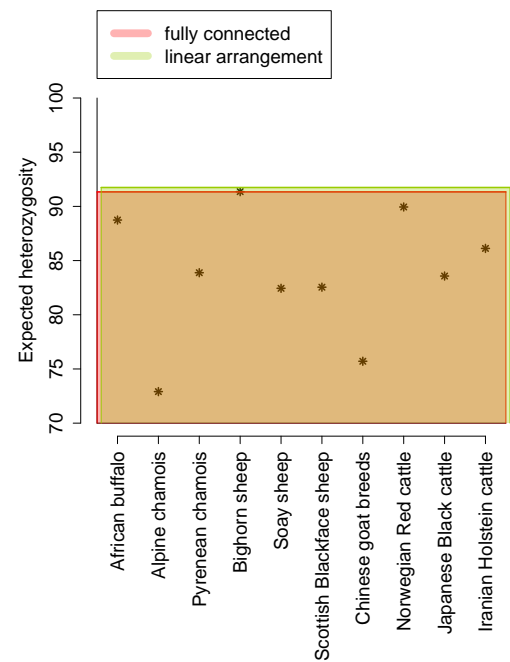
(a) Migration every generation



(b) Migration every 100 generations

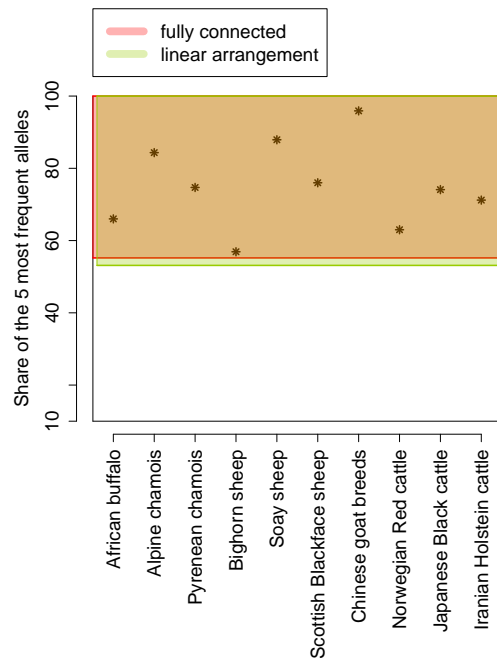


(c) Migration every 100000 generations

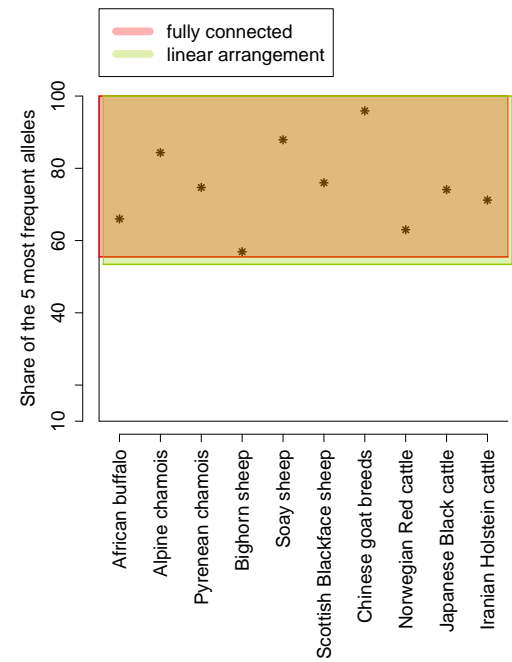


(d) No migration (isolated subpopulations)

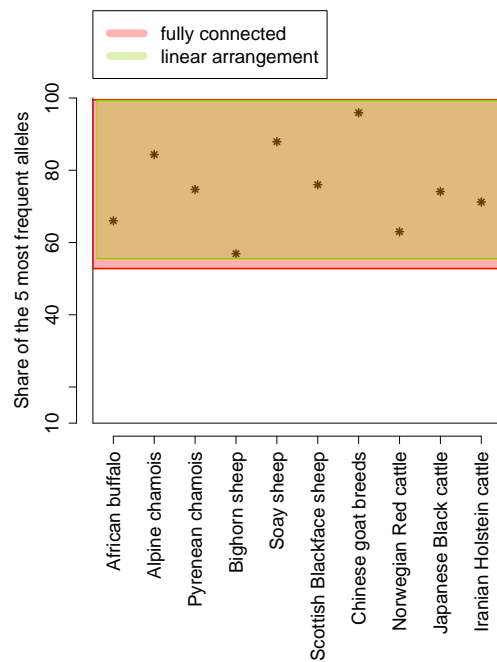
Figure D.28: Expected heterozygosity for metapopulations of 5 subpopulations with different connectivities and various migration rates, and for observed data. The red and green bands span population sizes from 10000 to 10 million. Figure D.28d gives an indication of the intrinsic stochastic variation, as the arrangement does not have an influence in a metapopulation of isolated subpopulations.



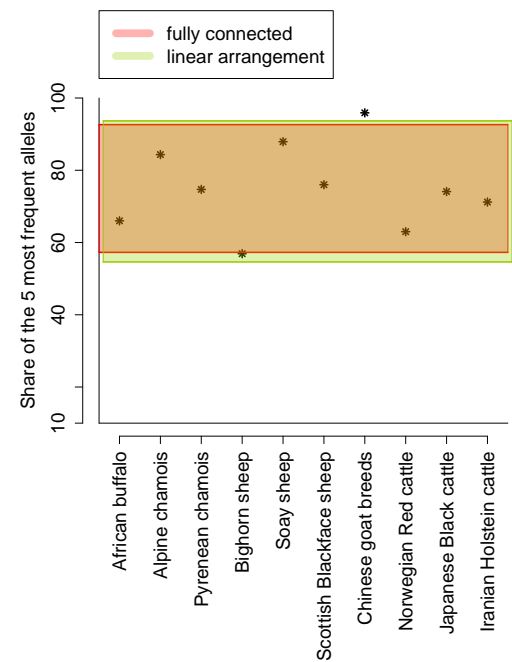
(a) Migration every generation



(b) Migration every 100 generations



(c) Migration every 100000 generations



(d) No migration (isolated subpopulations)

Figure D.29: Share of the five most frequent alleles for metapopulations of 5 subpopulations with different connectivities and various migration rates, and for observed data. The red and green bands span population sizes from 10000 to 10 million. Figure D.29d gives an indication of the intrinsic stochastic variation, as the arrangement does not have an influence in a metapopulation of isolated subpopulations.

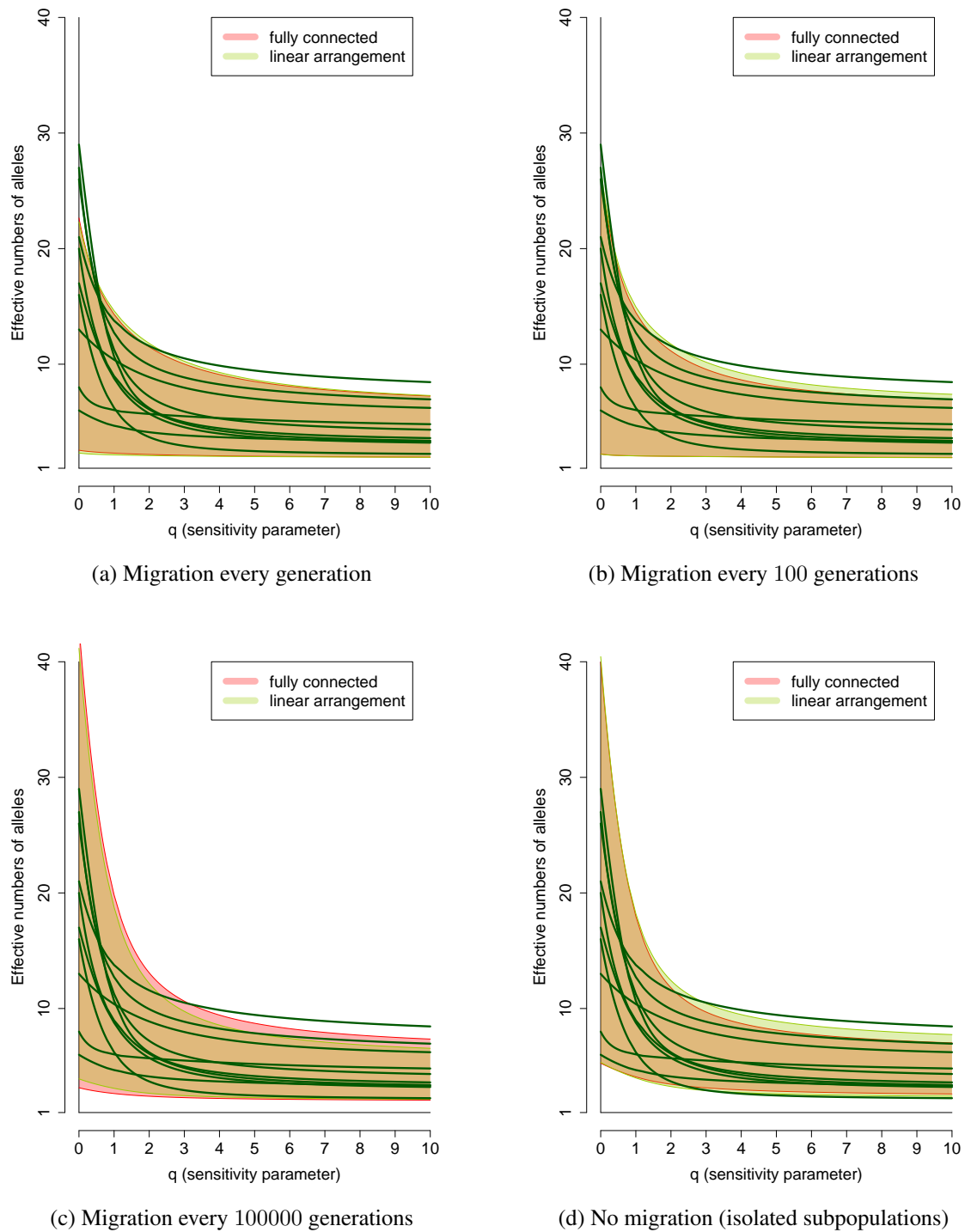
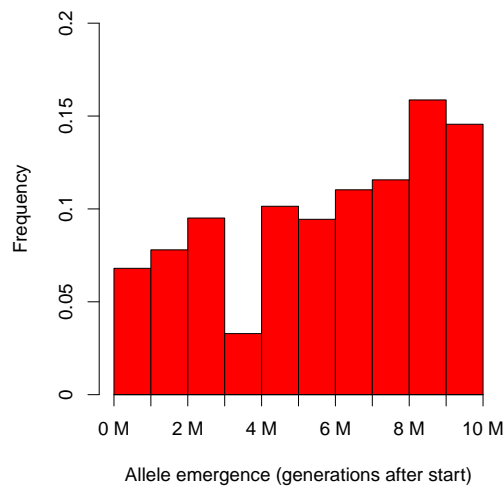
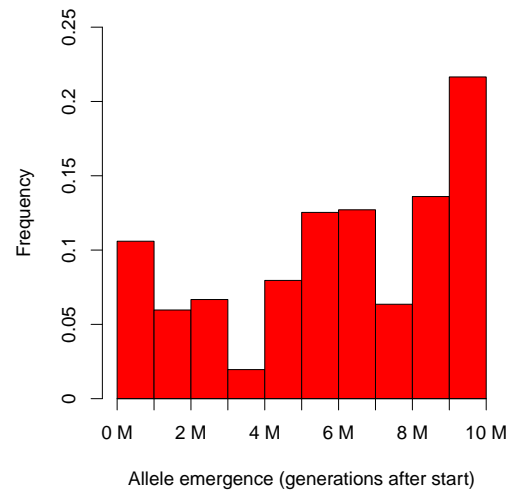


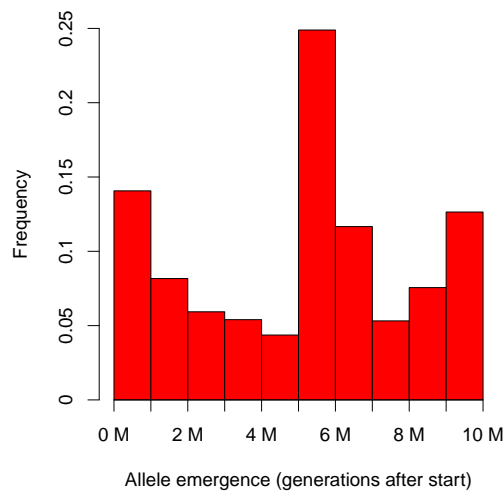
Figure D.30: Diversity profiles (naive diversity) for metapopulations with different connectivities of subpopulations and observed data for a metapopulation of 5 subpopulations and various migration rates. The red and green bands span population sizes from 10000 to 10 million. Figure D.30d gives an indication of the intrinsic stochastic variation, as the arrangement does not have an influence in a metapopulation of isolated subpopulations.



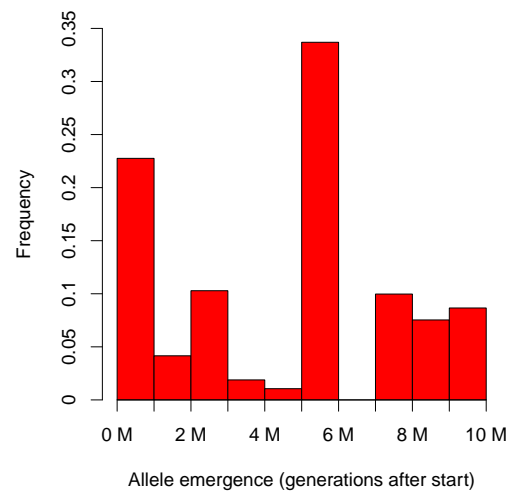
(a) 10000 ind. (total)



(b) 100000 ind. (total)

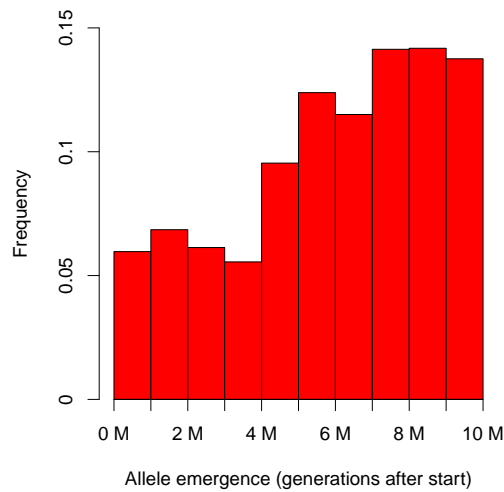


(c) 1 million ind. (total)

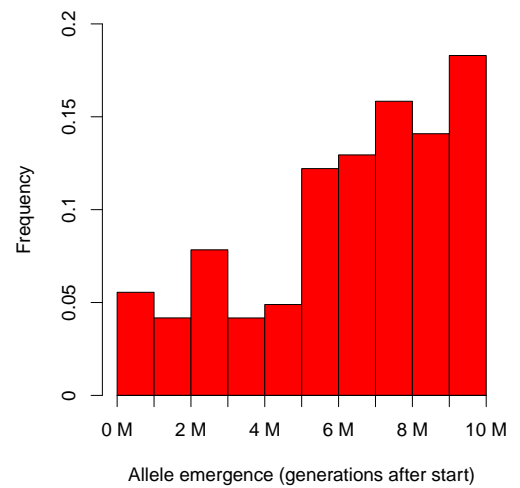


(d) 10 million ind. (total)

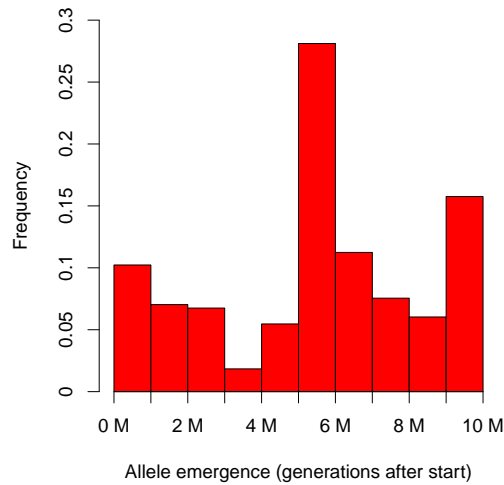
Figure D.31: Distribution of emergence times of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every generation.



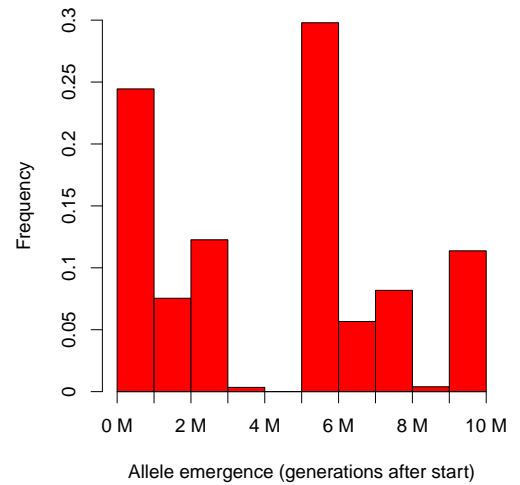
(a) 10000 ind. (total)



(b) 100000 ind. (total)



(c) 1 million ind. (total)



(d) 10 million ind. (total)

Figure D.32: Distribution of emergence times of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every 100 generations.

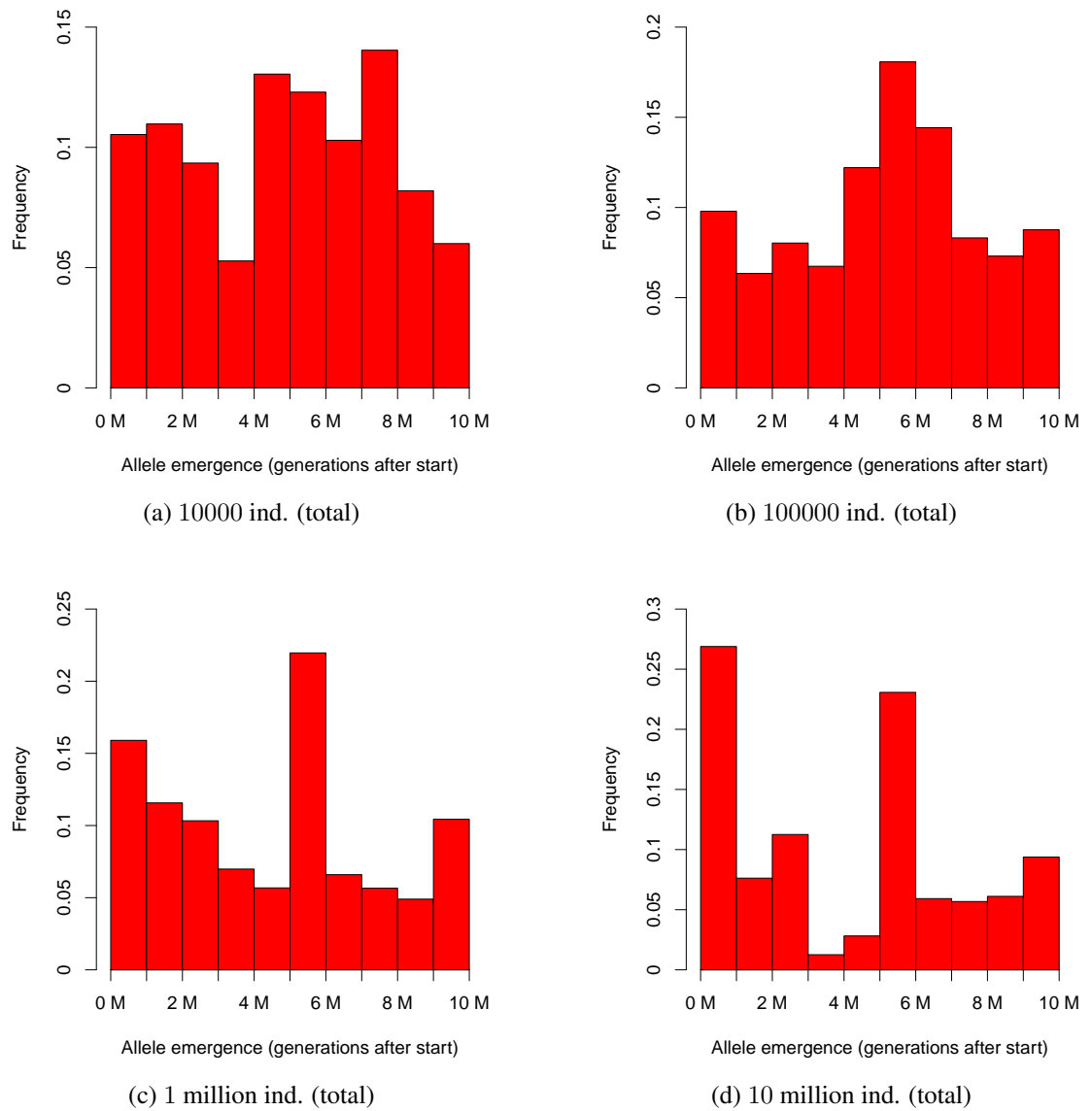
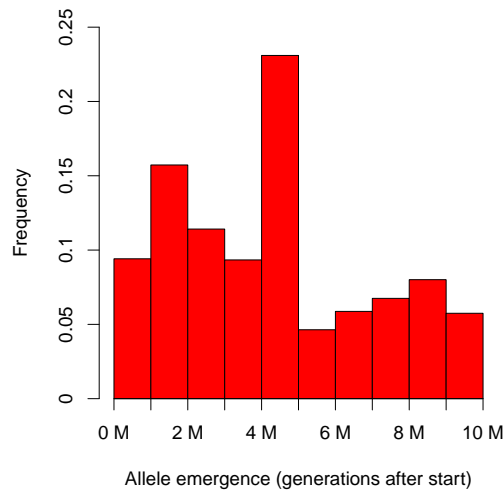
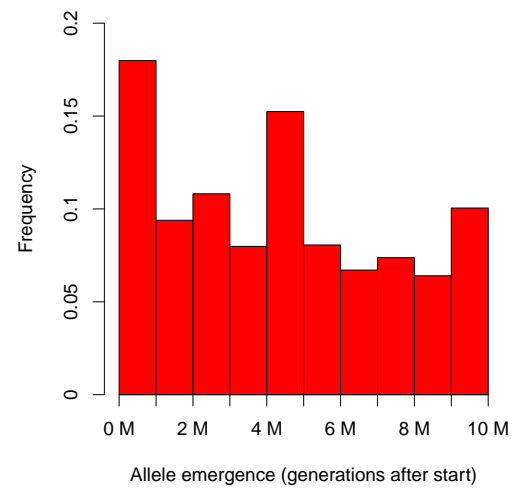


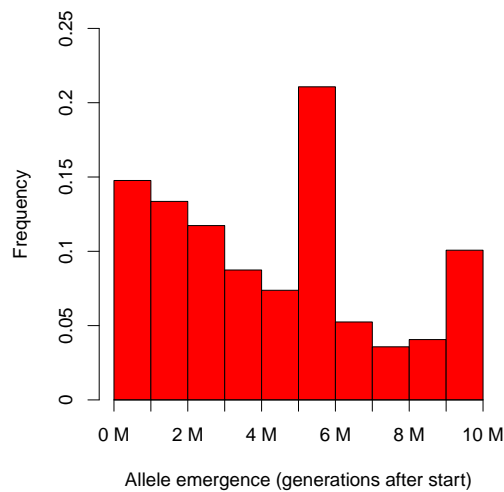
Figure D.33: Distribution of emergence times of alleles for a metapopulation in a star arrangement, all population sizes tested and a migration rate of 2 individuals every 100,000 generations.



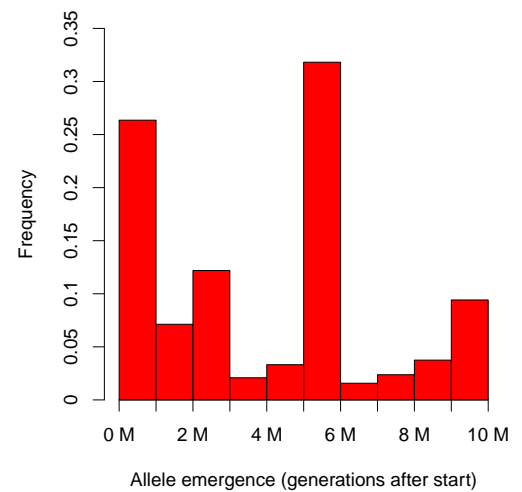
(a) 10000 ind. (total)



(b) 100000 ind. (total)

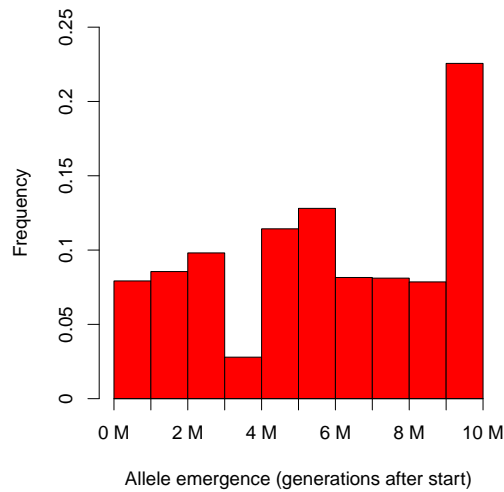


(c) 1 million ind. (total)

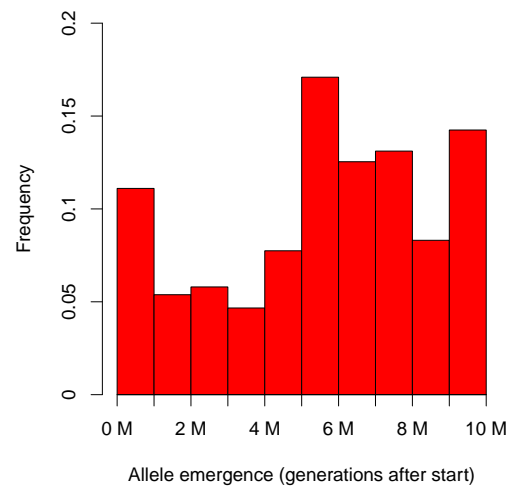


(d) 10 million ind. (total)

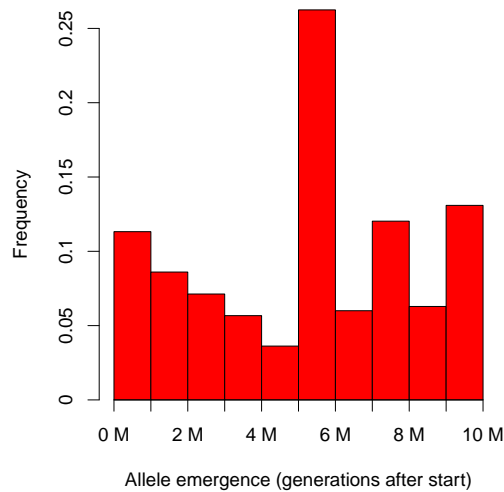
Figure D.34: Distribution of emergence times of alleles for a metapopulation in a star arrangement, all population sizes tested and a metapopulation of 5 isolated subpopulations.



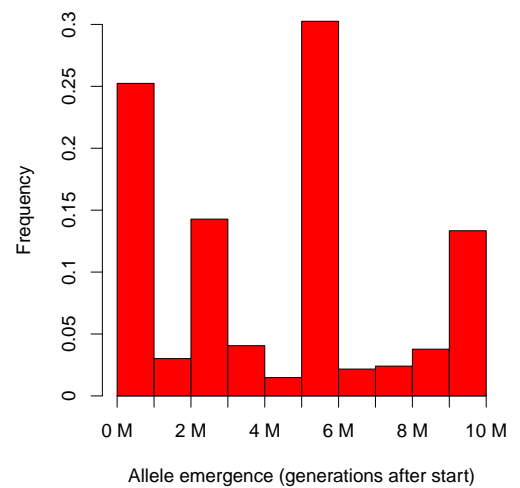
(a) 10000 ind. (total)



(b) 100000 ind. (total)

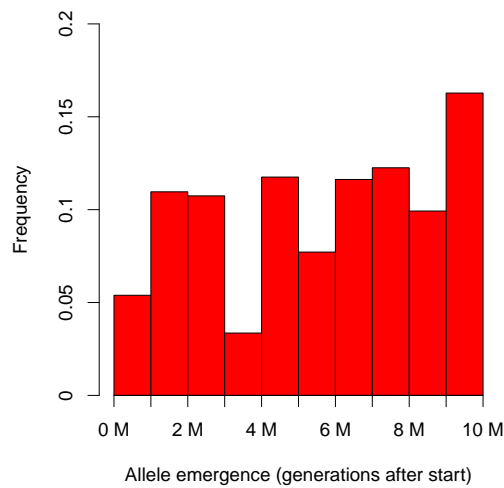


(c) 1 million ind. (total)

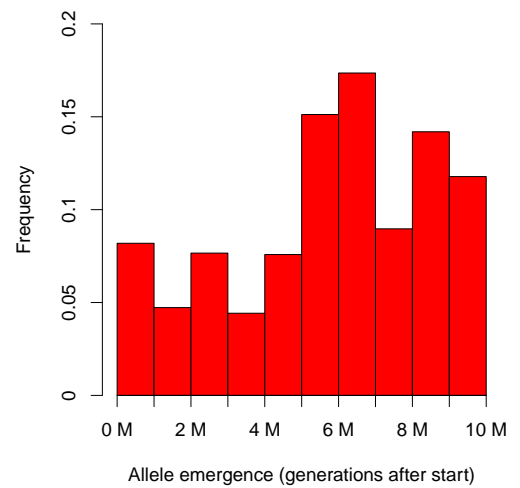


(d) 10 million ind. (total)

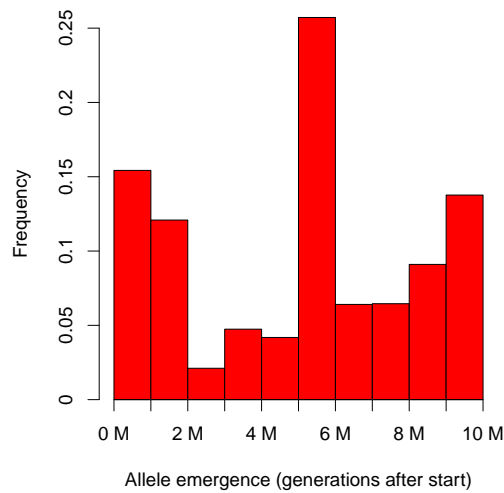
Figure D.35: Distribution of emergence times of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every generation.



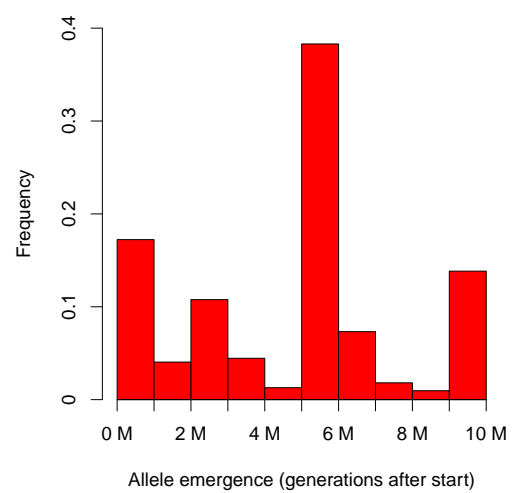
(a) 10000 ind. (total)



(b) 100000 ind. (total)

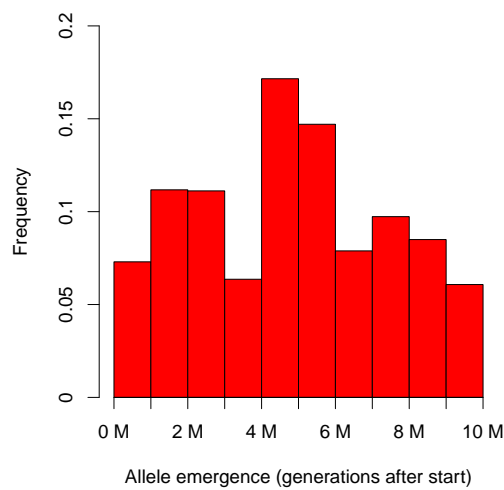


(c) 1 million ind. (total)

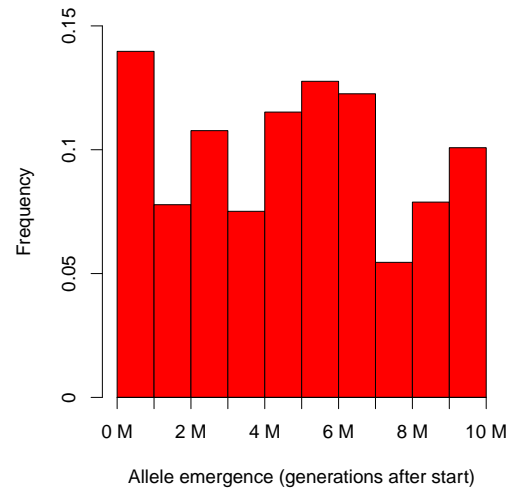


(d) 10 million ind. (total)

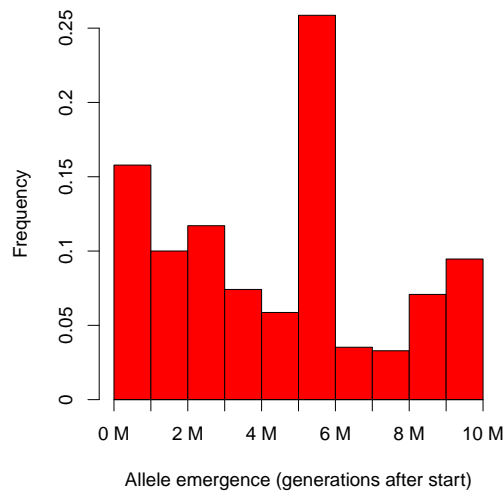
Figure D.36: Distribution of emergence times of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every 100 generations.



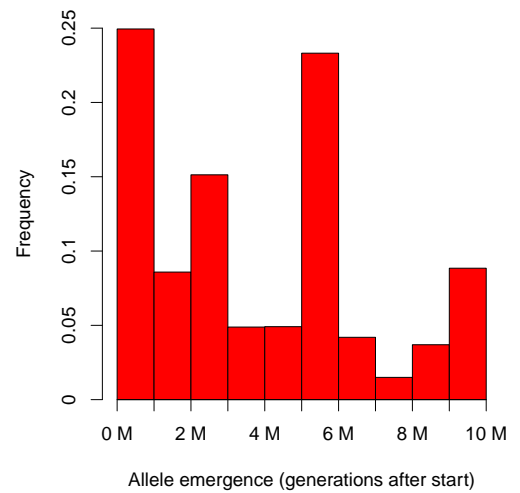
(a) 10000 ind. (total)



(b) 100000 ind. (total)

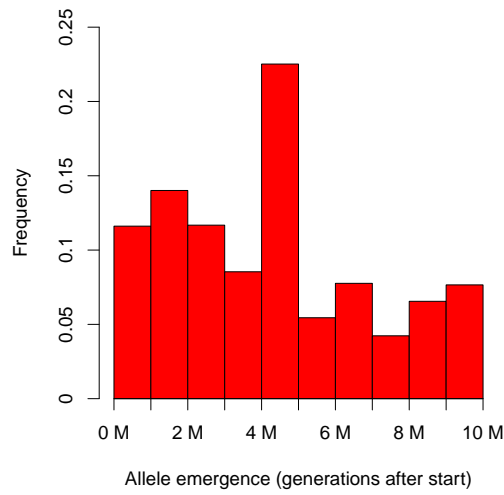


(c) 1 million ind. (total)

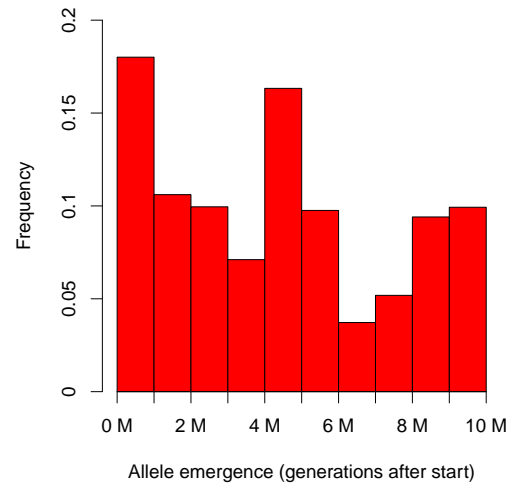


(d) 10 million ind. (total)

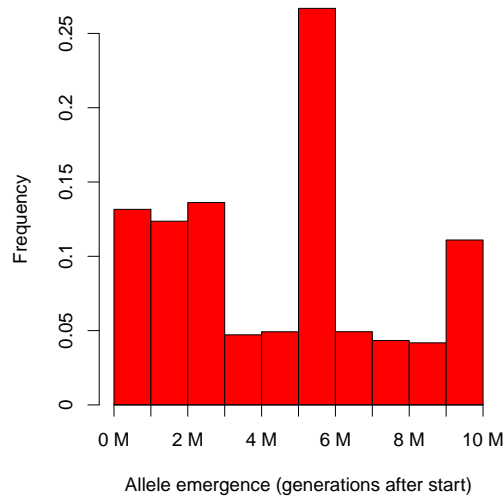
Figure D.37: Distribution of emergence times of alleles for a metapopulation in a linear arrangement, all population sizes tested and a migration rate of 2 individuals every 100000 generations.



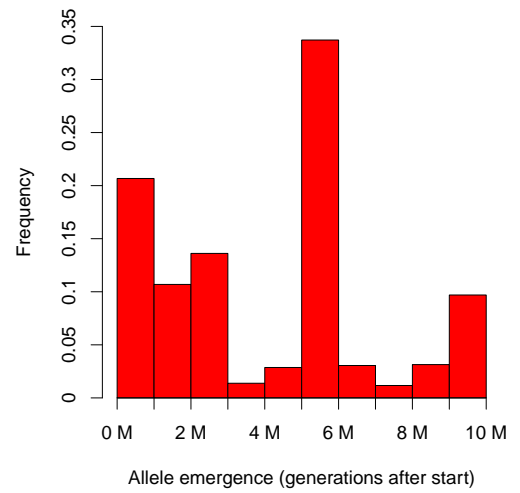
(a) 10000 ind. (total)



(b) 100000 ind. (total)



(c) 1 million ind. (total)



(d) 10 million ind. (total)

Figure D.38: Distribution of emergence times of alleles for a metapopulation in a linear arrangement, all population sizes tested and a metapopulation of 5 isolated subpopulations.

Bibliography

- Altuvia, Y. and Margalit, H. (2004). A structure-based approach for prediction of MHC-binding peptides. *Methods* 34, 454–459.
- Andersson, L., Lunden, A., Sigurdardottir, S., Davies, C. J. and Rask, L. (1988). Linkage relationships in the bovine MHC region. High recombination frequency between class II subregions. *Immunogenetics* 27, 273–280.
- Apanius, V., Penn, D., Slev, P. R., Ruff, R. L. and Potts, W. K. (1997). The Nature of Selection on the Major Histocompatibility Complex. *Critical Reviews in Immunology* 17, 179–224.
- Arkush, K. D., Giese, A. R., Mendonca, H. L., McBride, A. M., Marty, G. D. and Hedrick, P. W. (2002). Resistance to three pathogens in the endangered winter-run chinook salmon (*Oncorhynchus tshawytscha*): effects of inbreeding and major histocompatibility complex genotypes. *Canadian Journal of Fisheries and Aquatic Sciences* 59, 966–975.
- Ballingall, K. T., Rocchi, M. S., McKeever, D. J. and Wright, F. (2010). Trans-species polymorphism and selection in the MHC class II DRA genes of domestic sheep. *PloS ONE* 5, e11402.
- Berger, J. and Cunningham, C. (1995). Multiple bottlenecks, allopatric lineages and Badlands bison *Bos bison*: Consequences of lineage mixing. *Biological Conservation* 71, 13–23.
- Bernatchez, L. and Landry, C. (2003). MHC studies in nonmodel vertebrates: what have we learned about natural selection in 15 years? *Journal of Evolutionary Biology* 16, 363–377.
- Bjorkman, P. J., Saper, M. A., Samraoui, B., Bennett, W. S., Strominger, J. L. and Wiley, D. C. (1987). The foreign antigen binding site and T cell recognition regions of class I histocompatibility antigens. *Nature* 329, 512–518.

- Black, F. L. and Hedrick, P. W. (1997). Strong balancing selection at HLA loci: evidence from segregation in South Amerindian families. *Proceedings of the National Academy of Sciences of the United States of America* 94, 12452–12456.
- Bonneaud, C., Chastel, O., Federici, P., Westerdahl, H. and Sorci, G. (2006). Complex Mhc-based mate choice in a wild passerine. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 273, 1111–1116.
- Borghans, J. A. M., Beltman, J. B. and De Boer, R. J. (2004). MHC polymorphism under host-pathogen coevolution. *Immunogenetics* 55, 732–739.
- Bronson, P. G., Mack, S. J., Erlich, H. A. and Slatkin, M. (2013). A sequence-based approach demonstrates that balancing selection in classical human leukocyte antigen (HLA) loci is asymmetric. *Human Molecular Genetics* 22, 252–261.
- Brown, J. H., Jardetzky, T. S., Gorga, J. C., Stern, L. J., Urban, R. G., Strominger, J. L. and Wiley, D. C. (1993). Three-dimensional structure of the human class II histocompatibility antigen HLA-DR1. *Nature* 364, 33–39.
- Browning, M. and McMichael, A. (1996). *HLA and MHC: Genes, Molecules and Function*. Bios Scientific Publishers, Oxford.
- Carrington, M., Nelson, G. W., Martin, M. P., Kissner, T., Vlahov, D., Goedert, J. J., Kaslow, R., Buchbinder, S., Hoots, K. and O'Brien, S. J. (1999). HLA and HIV-1: Heterozygote Advantage and B*35-Cw*04 Disadvantage. *Science* 283, 1748–1752.
- Cavallero, S., Marco, I., Lavín, S., D'Amelio, S. and López-Olvera, J. R. (2012). Polymorphisms at MHC class II DRB1 exon 2 locus in Pyrenean chamois (*Rupicapra pyrenaica pyrenaica*). *Infection, Genetics and Evolution* 12, 1020–1026.
- Consuegra, S. and Garcia de Leaniz, C. (2008). MHC-mediated mate choice increases parasite resistance in salmon. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 275, 1397–1403.
- Crow, J. F. and Kimura, M. (1970). *An introduction to population genetics theory*. Harper & Row, New York.
- De Boer, R. J., Borghans, J. A. M., van Boven, M., Kesmir, C. and Weissing, F. J. (2004). Heterozygote advantage fails to explain the high degree of polymorphism of the MHC. *Immunogenetics* 55, 725–731.
- de Roos, A. P. W., Hayes, B. J., Spelman, R. J. and Goddard, M. E. (2008). Linkage Disequilibrium and Persistence of Phase in Holstein-Friesian, Jersey and Angus Cattle. *Genetics* 179, 1503–1512.

- Decker, J. E., Pires, J. C., Conant, G. C., McKay, S. D., Heaton, M. P., Chen, K., Cooper, A., Vilkki, J., Seabury, C. M., Caetano, A. R., Johnson, G. S., Brenneman, R. A., Hanotte, O., Eggert, L. S., Wiener, P., Kim, J.-j., Suk, K., Sonstegard, T. S., Tassell, C. P. V., Neibergs, H. L., Mcewan, J. C., Brauning, R., Coutinho, L. L., Babar, M. E., Wilson, G. A., McClure, M. C., Rolf, M. M., Kim, J., Schnabel, R. D. and Taylor, J. F. (2009). Resolving the evolution of extant and extinct ruminants with high-throughput phylogenomics. *Proceedings of the National Academy of Sciences of the United States of America* 106, 18644–18649.
- Dionne, M., Miller, K. M., Dodson, J. J. and Bernatchez, L. (2009). MHC standing genetic variation and pathogen resistance in wild Atlantic salmon. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 364, 1555–1565.
- Doherty, P. C. and Zinkernagel, R. M. (1975). A biological role for the major histocompatibility antigens. *The Lancet* 305, 1406–1409.
- Drummond, A. J., Rambaut, A., Shapiro, B. and Pybus, O. G. (2005). Bayesian coalescent inference of past population dynamics from molecular sequences. *Molecular Biology and Evolution* 22, 1185–1192.
- Edwards, S. V. and Hedrick, P. W. (1998). Evolution and ecology of MHC molecules: from genomics to sexual selection. *Trends in Ecology & Evolution* 13, 305–11.
- Eizaguirre, C., Lenz, T. L., Kalbe, M. and Milinski, M. (2012). Divergent selection on locally adapted major histocompatibility complex immune genes experimentally proven in the field. *Ecology Letters* 15, 723–731.
- Ejsmond, M. J., Babik, W. and Radwan, J. (2010). MHC allele frequency distributions under parasite-driven selection: A simulation model. *BMC Evolutionary Biology* 10, 332.
- EMBL-EBI (2015). Immuno Polymorphism Database.
- Evans, M. L. and Neff, B. D. (2009). Major histocompatibility complex heterozygote advantage and widespread bacterial infections in populations of Chinook salmon (*Oncorhynchus tshawytscha*). *Molecular Ecology* 18, 4716–4729.
- Ewens, W. J. (1972). The Sampling Theory of Selectively Neutral Alleles. *Theoretical Population Biology* 112, 87–112.
- Figuerola, F., Eberhardt, G. and Klein, J. (1988). MHC polymorphism pre-dating speciation. *Nature* 335, 265–267.
- Gautier, M., Faraut, T., Moazami-Goudarzi, K., Navratil, V., Foglio, M., Grohs, C., Boland, A., Garnier, J.-G., Boichard, D., Lathrop, G. M., Gut, I. G. and Eggen, A. (2007). Genetic

- and haplotypic structure in 14 European and African cattle breeds. *Genetics* 177, 1059–1070.
- Gillespie, J. H. (1977). A General Model to Account for Enzyme Variation in Natural Populations. III. Multiple Alleles. *Evolution* 31, 85–90.
- Goddard, M. E. and Hayes, B. J. (2009). Mapping genes for complex traits in domestic animals and their use in breeding programmes. *Nature Reviews Genetics* 10, 381–391.
- Grossen, C., Keller, L., Biebach, I., The International Goat Genome Consortium and Croll, D. (2014). Introgression from Domestic Goat Generated Variation at the Major Histocompatibility Complex of Alpine Ibex. *PLoS Genetics* 10, e1004438.
- Gutierrez-Espeleta, G. A., Hedrick, P. W., Kalinowski, S. T., Garrigan, D. and Boyce, W. M. (2001). Is the decline of desert bighorn sheep from infectious disease the result of low MHC variation? *Heredity* 86, 439–450.
- Gyllensten, U. B. and Erlich, H. A. (1989). Ancient roots for polymorphism at the HLA-DQ alpha locus in primates. *Proceedings of the National Academy of Sciences of the United States of America* 86, 9986–9990.
- Hamilton, M. S. and Hellström, I. (1978). Selection for Histoincompatible Progeny in Mice. *Biology of Reproduction* 19, 267–270.
- Havlicek, J. and Roberts, S. C. (2009). MHC-correlated mate choice in humans: a review. *Psychoneuroendocrinology* 34, 497–512.
- Hedrick, P. W. (1994). Evolutionary Genetics of the Major Histocompatibility Complex. *The American Naturalist* 143, 945–964.
- Hedrick, P. W. (1999). Balancing selection and MHC. *Genetica* 104, 207–214.
- Hedrick, P. W. (2002). Pathogen Resistance and Genetic Variation at MHC Loci. *Evolution* 56, 1902–1908.
- Hedrick, P. W. and Black, F. L. (1997). HLA and mate selection: no evidence in South Amerindians. *American Journal of Human Genetics* 61, 505–511.
- Hedrick, P. W. and Thomson, G. (1983). Evidence for balancing selection at HLA. *Genetics* 104, 449–456.
- Hill, A. V. S., Allsopp, C. E. M., Kwiatkowski, D., Anstey, N. M., Twumasi, P., Rowe, P. A., Bennett, S., Brewster, D., McMichael, A. J. and Greenwood, B. M. (1991). Common West African HLA antigens are associated with protection from severe malaria. *Nature* 352, 595–600.

- Hill, M. O. (1973). Diversity and Evenness: A Unifying Notation and Its Consequences. *Ecology* 54, 427–432.
- Hings, I. and Billingham, R. E. (1985). Maternal-fetal immune interactions and the maintenance of major histocompatibility complex polymorphism in the rat. *Journal of Reproductive Immunology* 7, 337–350.
- Hoekstra, R. F., Bijlsma, R. and Dolman, A. J. (1985). Polymorphism from environmental heterogeneity: models are only robust if the heterozygote is close in fitness to the favoured homozygote in each environment. *Genetical Research* 45, 299–314.
- Huchard, E., Knapp, L. A., Wang, J., Raymond, M. and Cowlshaw, G. (2010). MHC, mate choice and heterozygote advantage in a wild social primate. *Molecular Ecology* 19, 2545–2561.
- Hughes, A. L. and Nei, M. (1988). Patterns of nucleotide substitution at major histocompatibility complex class 1 loci reveals overdominant selection. *Nature* 335, 167–170.
- Hughes, A. L. and Nei, M. (1992). Maintenance of MHC Polymorphism. *Nature* 355, 402–403.
- Hughes, A. L. and Yeager, M. (1998). Natural selection at major histocompatibility complex loci of vertebrates. *Annual Review of Genetics* 32, 415–435.
- Ilmonen, P., Penn, D. J., Damjanovich, K., Morrison, L., Ghotbi, L. and Potts, W. K. (2007). Major histocompatibility complex heterozygosity reduces fitness in experimentally infected mice. *Genetics* 176, 2501–2508.
- Jeffery, K. J. M. and Bangham, C. R. M. (2000). Do infectious diseases drive MHC diversity? *Microbes and Infection* 2, 1335–1341.
- Jermann, T. M., Opitz, J. G., Stackhouse, J. and Benner, S. A. (1995). Reconstructing the evolutionary history of the artiodactyl ribonuclease superfamily. *Nature* 374, 57 – 59.
- Karlin, S. and Lessard, S. (1984). On the optimal sex-ratio: A stability analysis based on a characterization for one-locus multiallele viability models. *Journal of Mathematical Biology* 20, 15–38.
- Kekäläinen, J., Vallunen, J. A., Primmer, C. R., Rättä, J. and Taskinen, J. (2009). Signals of major histocompatibility complex overdominance in a wild salmonid population. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 276, 3133–3140.
- Kimura, M. (1983). *The Neutral Theory of Molecular Evolution*. Cambridge University Press, Cambridge.

- Kimura, M. and Crow, J. F. (1964). The number of alleles that can be maintained in a finite population. *Genetics* 49, 725–738.
- Klein, J. (1986). *Natural History of the Major Histocompatibility Complex*. John Wiley and Sons, New York.
- Klein, J., Satta, Y. and O'hUigin, C. (1993). The molecular descent of the major histocompatibility complex. *Annual Review of Immunology* 11, 269–295.
- Kulberg, S., Heringstad, B., Guttersrud, O. A. and Olsaker, I. (2007). Study on the association of BoLA-DRB3.2 alleles with clinical mastitis in Norwegian Red cows. *Journal of Animal Breeding and Genetics* 124, 201–207.
- Kumar, S. and Subramanian, S. (2002). Mutation rates in mammalian genomes. *Proceedings of the National Academy of Sciences of the United States of America* 99, 803–808.
- Lawlor, D. A., Ward, F. E., Ennis, P. D., Jackson, A. P. and Parham, P. (1988). HLA-A and B polymorphisms predate the divergence of humans and chimpanzees. *Nature* 335, 268–271.
- Leinster, T. and Cobbold, C. A. (2012). Measuring diversity: the importance of species similarity. *Ecology* 93, 477–489.
- Lenz, T. L. (2011). Computational prediction of MHC II-antigen binding supports divergent allele advantage and explains trans-species polymorphism. *Evolution* 65, 2380–2390.
- Lenz, T. L., Mueller, B., Trillmich, F. and Wolf, J. B. W. (2013). Divergent allele advantage at MHC-DRB through direct and maternal genotypic effects and its consequences for allele pool composition and mating. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 280, 1–9.
- Lenz, T. L., Wells, K., Pfeiffer, M. and Sommer, S. (2009). Diverse MHC IIB allele repertoire increases parasite resistance and body condition in the Long-tailed giant rat (*Leopoldamys sabanus*). *BMC Evolutionary Biology* 9, 269.
- Lewin, H. A., Russell, G. C. and Glass, E. J. (1999). Comparative organization and function of the major histocompatibility complex of domesticated cattle. *Immunological Reviews* 167, 145–158.
- Lewontin, R. C., Ginzburg, L. R. and Tuljapurkar, S. D. (1978). Heterosis as an Explanation for Large Amounts of Genic Polymorphism. *Genetics* 88, 149–169.
- Li, M.-H., Li, K., Kantanen, J., Feng, Z., Fan, B. and Zhao, S.-H. (2006). Allelic variations in exon 2 of caprine MHC class II DRB3 gene in Chinese indigenous goats. *Small Ruminant Research* 66, 236–243.

- Li, W.-H. (1978). Maintenance of Genetic Variability under the Joint Effect of Mutation, Selection and Random Drift. *Genetics* 90, 349–382.
- Liberman, U. (1991). On the Relation Between the Instability of ESS in Discrete Dynamics and Segregation Distortion. *Journal of Theoretical Biology* 150, 421–436.
- MacEachern, S., Hayes, B., McEwan, J. and Goddard, M. (2009). An examination of positive selection and changing effective population size in Angus and Holstein cattle populations using a high density SNP genotyping platform and the contribution of ancient polymorphism to genomic diversity in domestic cattle. *BMC Genomics* 10, 181.
- Marks, R. W. and Spencer, H. G. (1991). The Maintenance of Single-Locus Polymorphism. II. The Evolution of Fitnesses and Allele Frequencies. *The American Naturalist* 138, 1354–1371.
- Maruyama, T. and Nei, M. (1981). Genetic variability maintained by mutation and overdominant selection in finite populations. *Genetics* 98, 441–459.
- Mayer, W. E., Jonker, M., Klein, D., Ivanyi, P., van Seventer, G. and Klein, J. (1988). Nucleotide sequences of chimpanzee MHC class I alleles: evidence for trans-species mode of evolution. *The EMBO Journal* 7, 2765–2774.
- McClelland, E. E., Penn, D. J. and Potts, W. K. (2003). Major histocompatibility complex heterozygote superiority during coinfection. *Infection and Immunity* 71, 2079–2086.
- Menotti-Raymond, M. and O'Brien, S. J. (1993). Dating the genetic bottleneck of the African cheetah. *Proceedings of the National Academy of Sciences of the United States of America* 90, 3172–3176.
- Meyer, D. and Thomson, G. (2001). How selection shapes variation of the human major histocompatibility complex: a review. *Annals of Human Genetics* 65, 1–26.
- Meyer-Lucht, Y. and Sommer, S. (2005). MHC diversity and the association to nematode parasitism in the yellow-necked mouse (*Apodemus flavicollis*). *Molecular Ecology* 14, 2233–2243.
- Mikko, S., Røed, K., Schmutz, S. and Andersson, L. (1999). Monomorphism and polymorphism at Mhc DRB loci in domestic and wild ruminants. *Immunological reviews* 167, 169–178.
- Milinski, M. (2006). The Major Histocompatibility Complex, Sexual Selection, and Mate Choice. *Annual Review of Ecology, Evolution, and Systematics* 37, 159–186.

- Miyasaka, T., Takeshima, S.-n., Matsumoto, Y., Kobayashi, N., Matsuhashi, T., Miyazaki, Y., Tanabe, Y., Ishibashi, K., Sentsui, H. and Aida, Y. (2011). The diversity of bovine MHC class II DRB3 and DQA1 alleles in different herds of Japanese Black and Holstein cattle in Japan. *Gene* 472, 42–49.
- Nachman, M. W. and Crowell, S. L. (2000). Estimate of the mutation rate per nucleotide in humans. *Genetics* 156, 297–304.
- Nassiry, M. R., Eftekhari Shahroodi, F., Mosafer, J., Mohammadi, A., Manshad, E., Ghazanfari, S., Mohammad Abadi, M. R. and Sulimova, G. E. (2005). Analysis and frequency of bovine lymphocyte antigen (BoLA-DRB3) alleles in Iranian Holstein cattle. *Genetika* 41, 817–822.
- Nei, M. (1987). *Molecular Evolutionary Genetics*. Columbia University Press, New York.
- O'Connor, S. L., Lhost, J. J., Becker, E. A., Detmer, A. M., Johnson, R. C., Macnair, C. E., Wiseman, R. W., Karl, J. A., Greene, J. M., Burwitz, B. J., Bimber, B. N., Lank, S. M., Tuscher, J. J., Mee, E. T., Rose, N. J., Desrosiers, R. C., Hughes, A. L., Friedrich, T. C., Carrington, M. and O'Connor, D. H. (2010). MHC heterozygote advantage in simian immunodeficiency virus-infected Mauritian cynomolgus macaques. *Science Translational Medicine* 2, 22ra18.
- Ohashi, J., Naka, I., Toyoda, A., Takasu, M., Tokunaga, K., Ishida, T., Sakaki, Y. and Hohjoh, H. (2006). Estimation of the species-specific mutation rates at the DRB1 locus in humans and chimpanzee. *Tissue Antigens* 68, 427–31.
- Ohta, T. (1995). Gene conversion vs point mutation in generating variability at the antigen recognition site of major histocompatibility complex loci. *Journal of Molecular Evolution* 41, 115–119.
- Oliver, M. K., Telfer, S. and Pierniey, S. B. (2009). Major histocompatibility complex (MHC) heterozygote superiority to natural multi-parasite infections in the water vole (*Arvicola terrestris*). *Proceedings of the Royal Society of London. Series B, Biological Sciences* 276, 1119–1128.
- Parham, P. and Ohta, T. (1996). *Population Biology of Antigen Presentation by MHC Class I Molecules*. *Science* 272, 67–74.
- Paterson, S., Wilson, K. and Pemberton, J. M. (1998). Major histocompatibility complex variation associated with juvenile survival and parasite resistance in a large unmanaged ungulate population (*Ovis aries* L.). *Proceedings of the National Academy of Sciences of the United States of America* 95, 3714–3719.

- Penn, D. J. (2002). The Scent of Genetic Compatibility: Sexual Selection and the Major Histocompatibility Complex. *Ethology* 108, 1–21.
- Penn, D. J., Damjanovich, K. and Potts, W. K. (2002). MHC heterozygosity confers a selective advantage against multiple-strain infections. *Proceedings of the National Academy of Sciences of the United States of America* 99, 11260–11264.
- Piertney, S. B. and Oliver, M. K. (2006). The evolutionary ecology of the major histocompatibility complex. *Heredity* 96, 7–21.
- Potts, W. K., Manning, C. J. and Wakeland, E. K. (1994). The role of infectious disease, inbreeding and mating preferences in maintaining MHC genetic diversity: an experimental test. *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 346, 369–378.
- Potts, W. K. and Slev, P. R. (1995). Pathogen-based models favoring MHC genetic diversity. *Immunological Reviews* 143, 181–197.
- Potts, W. K. and Wakeland, E. K. (1990). Evolution of Diversity at the Major Histocompatibility Complex. *Trends in Ecology & Evolution* 5, 181–187.
- R Development Core Team (2011). R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing Vienna, Austria. ISBN 3-900051-07-0.
- Radwan, J., Biedrzycka, A. and Babik, W. (2010). Does reduced MHC diversity decrease viability of vertebrate populations? *Biological Conservation* 143, 537–544.
- Rammensee, H.-G., Friede, T. and Stevanovic, S. (1995). MHC ligands and peptide motifs: first listing. *Immunogenetics* 41, 178–228.
- Richman, A. D., Herrera, L. G. and Nash, D. (2001). MHC class II beta sequence diversity in the deer mouse (*Peromyscus maniculatus*): implications for models of balancing selection. *Molecular Ecology* 10, 2765–2773.
- Roach, J. C., Glusman, G., Smit, A. F. A., Huff, C. D., Hubley, R., Shannon, P. T., Rowen, L., Pant, K. P., Goodman, N., Bamshad, M., Shendure, J., Drmanac, R., Jorde, L. B., Hood, L. and Galas, D. J. (2010). Analysis of genetic inheritance in a family quartet by whole-genome sequencing. *Science* 328, 636–639.
- Robertson, A. (1962). Selection for Heterozygotes in Small Populations. *Genetics* 47, 1291–1300.
- Satta, Y. (1997). Effects of intra-locus recombination on HLA polymorphism. *Hereditas* 127, 105–112.

- Satta, Y., O'hUigin, C., Takahata, N. and Klein, J. (1994). Intensity of natural selection at the major histocompatibility complex loci. *Proceedings of the National Academy of Sciences of the United States of America* 91, 7184–7188.
- Sauermann, U., Nürnberg, P., Bercovitch, F. B., Berard, J. D., Trefilov, A., Widdig, A., Kessler, M., Schmidtke, J. and Krawczak, M. (2001). Increased reproductive success of MHC class II heterozygous males among free-ranging rhesus macaques. *Human Genetics* 108, 249–254.
- Schad, J., Ganzhorn, J. U. and Sommer, S. (2005). Parasite burden and constitution of major histocompatibility complex in the Malagasy mouse lemur, *Microcebus murinus*. *Evolution* 59, 439–50.
- Schaschl, H., Suchentrunk, F., Morris, D. L., Ben Slimen, H., Smith, S. and Arnold, W. (2012). Sex-specific selection for MHC variability in Alpine chamois. *BMC Evolutionary Biology* 12, 20.
- Schwaiger, F. W., Gostomski, D., Stear, M. J., Duncan, J. L., McKellar, Q. A., Epplen, J. T. and Buitkamp, J. (1995). An ovine major histocompatibility complex DRB1 allele is associated with low faecal egg counts following natural, predominantly *Ostertagia circumcincta* infection. *International Journal for Parasitology* 25, 815–822.
- Sena, L., Schneider, M. P. C., Brenig, B. B., Honeycutt, R. L., Honeycutt, D. A., Womack, J. E. and Skow, L. C. (2011). Polymorphism and gene organization of water buffalo MHC-DQB genes show homology to the BoLA DQB region. *Animal Genetics* 42, 378–385.
- Shapiro, B., Drummond, A. J., Rambaut, A., Wilson, M. C., Matheus, P. E., Sher, A. V., Pybus, O. G., Gilbert, M. T. P., Barnes, I., Binladen, J., Willerslev, E., Hansen, A. J., Baryshnikov, G. F., Burns, J. A., Davydov, S., Driver, J. C., Froese, D. G., Harington, C. R., Keddle, G., Kosintsev, P., Kunz, M. L., Martin, L. D., Stephenson, R. O., Storer, J., Tedford, R., Zimov, S. and Cooper, A. (2004). Rise and fall of the Beringian steppe bison. *Science* 306, 1561–1565.
- Slade, R. W. and McCallum, H. I. (1992). Overdominant vs. Frequency-Dependent Selection at MHC Loci. *Genetics* 132, 861–862.
- Sommer, S. (2005). The importance of immune gene variability (MHC) in evolutionary ecology and conservation. *Frontiers in Zoology* 2, 16.
- Spencer, H. G. and Marks, R. W. (1988). The Maintenance of Single-Locus Polymorphism. I. Numerical Studies of a Viability Selection Model. *Genetics* 120, 605–613.

- Spencer, H. G. and Marks, R. W. (1992). The Maintenance of Single-Locus Polymorphism. IV. Models With Mutation From Existing Alleles. *Genetics* 130, 211–221.
- Spurgin, L. G. and Richardson, D. S. (2010). How pathogens drive genetic diversity: MHC, mechanisms and misunderstandings. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 277, 979–988.
- Stear, M. J., Innocent, G. T. and Buitkamp, J. (2005). The evolution and maintenance of polymorphism in the major histocompatibility complex. *Veterinary Immunology and Immunopathology* 108, 53–57.
- Stoffels, R. J. and Spencer, H. G. (2008). An asymmetric model of heterozygote advantage at major histocompatibility complex genes: degenerate pathogen recognition and intersection advantage. *Genetics* 178, 1473–1489.
- Strachan, T. and Read, A. (2010). *Human Molecular Genetics*. 4th ed. edition, Garland Science, New York; Abingdon, Oxon.
- Takahata, N. and Nei, M. (1990). Allelic genealogy under overdominant and frequency-dependent selection and polymorphism of major histocompatibility complex loci. *Genetics* 124, 967–978.
- Takahata, N., Satta, Y. and Klein, J. (1992). Polymorphism and balancing selection at major histocompatibility complex loci. *Genetics* 130, 925–938.
- Thursz, M. R., Thomas, H. C., Greenwood, B. M. and Hill, A. V. (1997). Heterozygote advantage for HLA class-II type in hepatitis B virus infection. *Nature Genetics* 17, 11–12.
- Trowsdale, J. (1993). Genomic structure and function in the MHC. *Trends in Genetics* 9, 117–122.
- van Oosterhout, C. (2009). A new theory of MHC evolution: beyond selection on the immune genes. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 276, 657–665.
- Wakeland, E. K., Boehme, S., She, J. X., Lu, C. C., McIndoe, R. A., Cheng, I., Ye, Y. and Potts, W. K. (1990). Ancestral polymorphisms of MHC class II genes: divergent allele advantage. *Immunologic Research* 9, 115–122.
- Watterson, G. A. (1978). The homozygosity test of neutrality. *Genetics* 88, 405–417.
- Wedekind, C., Seebeck, T., Bettens, F. and Paepke, A. J. (1995). MHC-dependent mate preferences in humans. *Proceedings of the Royal Society of London. Series B, Biological Sciences* 260, 245–249.

- Weissing, F. J. and van Boven, M. (2001). Selection and segregation distortion in a sex-differentiated population. *Theoretical Population Biology* 60, 327–41.
- Wenink, P. W., Groen, A. F., Roelke-Parker, M. E. and Prins, H. H. T. (1998). African buffalo maintain high genetic diversity in the major histocompatibility complex in spite of historically known population bottlenecks. *Molecular Ecology* 7, 1315–1322.
- Worley, K., Collet, J., Spurgin, L. G., Cornwallis, C., Pizzari, T. and Richardson, D. S. (2010). MHC heterozygosity and survival in red junglefowl. *Molecular Ecology* 19, 3064–3075.
- Wright, S. (1990). Evolution in Mendelian populations. *Bulletin of Mathematical Biology* 52, 241–295.
- Wright, S. and Dobzhansky, T. (1946). Genetics of natural populations. XII. Experimental reproduction of some of the changes caused by natural selection in certain populations of *Drosophila pseudoobscura*. *Genetics* 31, 125–156.
- Wu, C. I. and Li, W. H. (1985). Evidence for higher rates of nucleotide substitution in rodents than in man. *Proceedings of the National Academy of Sciences of the United States of America* 82, 1741–1745.
- Yamazaki, B. Y. K., Boyse, E. A., Mike, V., Thaler, H. T., Mathieson, B. J., Abbott, J., Boyse, J., Zayas, Z. A. and Thomas, L. (1976). Control of mating preferences in mice by genes in the Major Histocompatibility Complex. *Journal of Experimental Medicine* 144, 1324–1335.
- Yokoyama, S. and Nei, M. (1979). Population dynamics of sex-determining alleles in honey bees and self-incompatibility alleles in plants. *Genetics* 91, 609–626.
- Yuhki, N. and O'Brien, S. J. (1990). DNA variation of the mammalian major histocompatibility complex reflects genomic diversity and population history. *Proceedings of the National Academy of Sciences of the United States of America* 87, 836–840.